

# How does generalization behave under suitable model capacities in modern machine learning

- A new  $\varphi$ -curve under norm-based capacity
- 

Fanghui Liu

[fanghui.liu@warwick.ac.uk](mailto: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 Mathematics, The University of Hong Kong



# My research

## Research interests

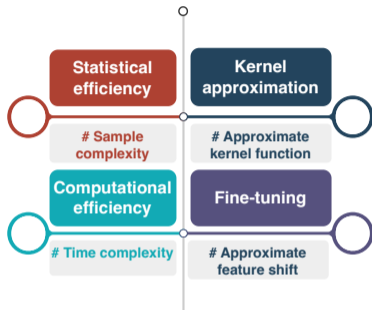
- Foundations of machine learning (ML)
- Theory-grounded efficient algorithm design
- Trustworthy ML



# My research

## Research interests

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



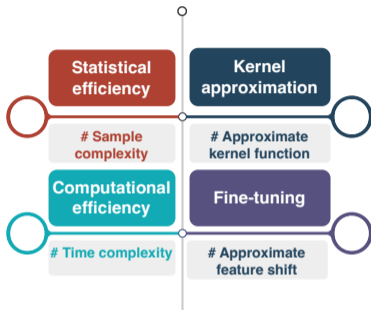
# My research

## Research interests

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

## Research goal

- characterize learning efficiency in theory
- contribute to practice



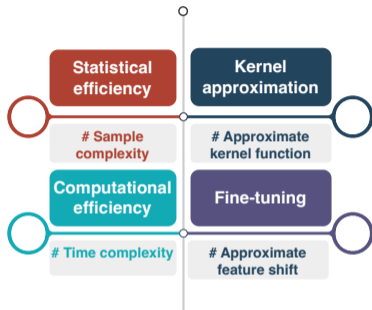
# My research

## 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)

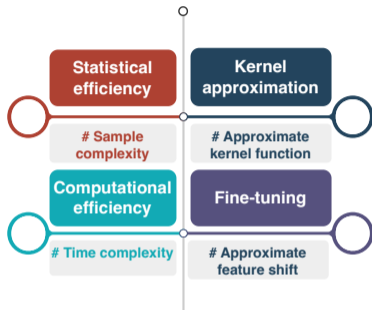
# My research

## 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)

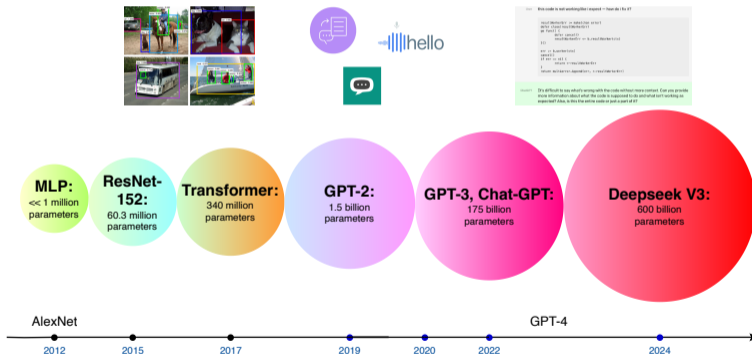
Machine learning works in high dimensions that can be a curse!

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

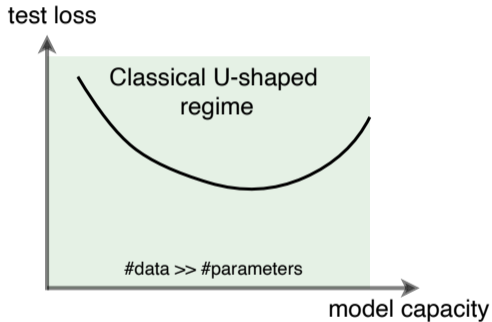


# In the era of machine learning

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

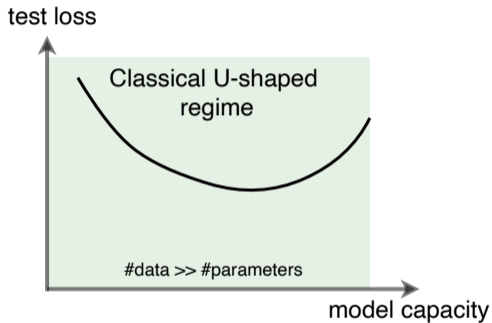


## ML textbooks: Larger models tend to overfit!

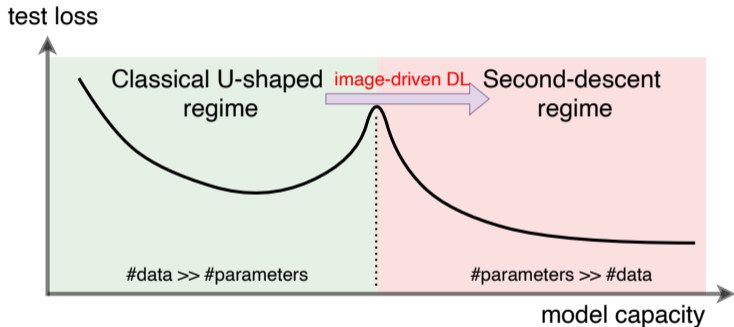


ML textbooks: Larger models tend to overfit!

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

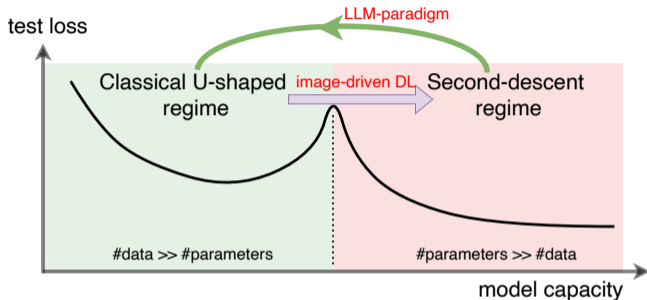


Figure 1: Paradigm among test loss, data, and model capacity.

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

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

# Learning paradigm in the past twenty years

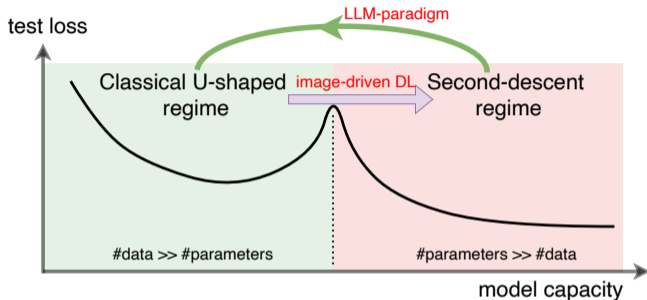


Figure 1: Paradigm among test loss, data, and model capacity.

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

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

# 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)

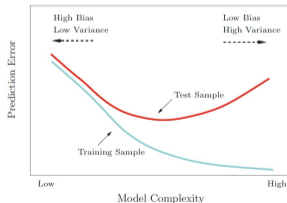
# A fundamental concept in machine learning: model capacity

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)

Trevor Hastie  
Robert Tibshirani  
Jerome Friedman

## The Elements of Statistical Learning

Data Mining, Inference, and Prediction

# 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!

# Today's talk: Norm-based capacity via deterministic equivalence

---

## Today's talk: Norm-based capacity via deterministic equivalence

(Bartlett et al., 2005)

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

(Bartlett, 1999)

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)

(Bartlett, 1999)

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: Norm-based capacity via deterministic equivalence

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?

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

- p How to precisely characterize the relationship under norm-based model capacity?
- ^ Reshape bias-variance trade-offs, double descent, scaling law under 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](#)

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

- p How to precisely characterize the relationship under norm-based model capacity?
  - ^ Reshape bias-variance trade-offs, double descent, scaling law under 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](#)
- p 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. IMLR 2024.

## Background: Random features ridge regression

$$f_p(x; \mathbf{a}) = \sum_{i=1}^p \mathbf{a}_i \phi_i(x; \mathbf{w}_i); \quad \phi_i := f(\mathbf{a}_i; \mathbf{w}_i) g_{i=1}^p$$

$\phi_i : X \rightarrow \mathbb{R}$ , e.g., ReLU:

$$\phi_i(x; \mathbf{w}_i) = \max(\mathbf{w}_i^T x; 0)$$

Random features models (RFMs) (Rahimi and Recht, 2007; Liu et al., 2021):

$\phi_i \mathbf{w}_i g_{i=1}^p \stackrel{\text{iid}}{\sim}$  for a given  $\mathbf{W} \in \mathbb{R}^{n \times p}$   
only train the second layer

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

$Z \in \mathbb{R}^{n \times p}$  with  $[Z]_{ij} = \phi_j(x_i; \mathbf{w}_j)$ .

Norm over the first-layer (untrained)  $\|\mathbf{w}_i\|_2^2$

Norm over the second-layer  $\|\mathbf{a}\|_2^2$

## Background: Random features ridge regression

$$f_p(x; \mathbf{a}) = \sum_{i=1}^p \mathbf{a}_i \phi(x; \mathbf{w}_i); \quad \phi(x; \mathbf{w}_i) := f(\mathbf{a}_i; \mathbf{w}_i) g_{i=1}^p$$

$\phi : X \times W \rightarrow \mathbb{R}$ , e.g., ReLU:

$$\phi(x; \mathbf{w}) = \max(\mathbf{w}^\top x; 0)$$

$\hat{\mathbf{a}}$  Random features models (RFMs) (Rahimi and Recht, 2007; Liu et al., 2021):

$\phi(\mathbf{w}_i) g_{i=1}^p$  iid for a given  $\mathbf{W} \in \mathbb{R}^{n \times p}$   
only train the second layer

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

$\mathbf{Z} \in \mathbb{R}^{n \times p}$  with  $[Z]_{ij} = \phi(x_i; \mathbf{w}_j)$ .

$\hat{\mathbf{a}}$  Norm over the first-layer (untrained)  $\|\mathbf{W}\|_F$

$\hat{\mathbf{a}}$  Norm over the second-layer  $\|\mathbf{a}\|_2^2$

## Background: Random features ridge regression

$$f_p(x; \mathbf{a}) = \sum_{i=1}^p \mathbf{a}_i \phi(x; \mathbf{w}_i); \quad \phi(x; \mathbf{w}_i) := f(\mathbf{a}_i; \mathbf{w}_i) g_{i=1}^p$$

$\phi : X \times W \rightarrow \mathbb{R}$ , e.g., ReLU:

$$\phi(x; \mathbf{w}) = \max(\langle \mathbf{x}; \mathbf{w}; 0)$$

Random features models (RFMs) (Rahimi and Recht, 2007; Liu et al., 2021):

$\phi(\mathbf{w}_i) g_{i=1}^p$  iid for a given  $\mathbf{W} \in \mathbb{R}^{n \times p}$   
only train the second layer

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

$$\mathbf{Z} \in \mathbb{R}^{n \times p} \text{ with } [\mathbf{Z}]_{ij} = \frac{1}{\sqrt{m}} \phi(\mathbf{x}_i; \mathbf{w}_j).$$

Norm over the first-layer (untrained)  $\|\mathbf{W}\|_F$

Norm over the second-layer  $\|\mathbf{a}\|_2^2$

## Background: Random features ridge regression

$$f_p(x; \mathbf{a}) = \sum_{i=1}^p \mathbf{a}_i \phi(x; \mathbf{w}_i); \quad \phi(x; \mathbf{w}_i) := f(\mathbf{a}_i; \mathbf{w}_i) g_{i=1}^p$$

$\phi : X \times \mathbb{R}^d \rightarrow \mathbb{R}$ , e.g., ReLU:

$$\phi(x; \mathbf{w}) = \max(\mathbf{w}^\top x; 0)$$

Random features models (RFMs) (Rahimi and Recht, 2007; Liu et al., 2021):

$\phi(\mathbf{w}_i) g_{i=1}^p$  iid for a given  $\mathbf{w} \in \mathbb{R}^d$   
only train the second layer

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

$Z \in \mathbb{R}^{n \times p}$  with  $[Z]_{ij} = \phi(x_i; \mathbf{w}_j)$ .

Norm over the first-layer (untrained)  $\|\mathbf{w}\|_F$

Norm over the second-layer  $\|\mathbf{a}\|_2^2$

## Background: Test risk of random features model

A compact integral operator  $T : L^2(\mathcal{X}) \rightarrow L^2(\mathcal{W})$  for any  $f \in L^2(\mathcal{X})$   
(De Lippis et al., 2024)

$$(Tf)(w) := \int_{\mathcal{R}^d} \kappa(x; w) f(x) d(x); \quad T = \sum_{k=1}^K \kappa_k \kappa_k'$$

Covariate feature matrix  $G := [g_1; \dots; g_n] \in \mathbb{R}^{n \times 1}$  with  $g_i := (\kappa_k(x_i))_{k=1}^K$

Weight feature matrix  $H := [h_1; \dots; h_p] \in \mathbb{R}^{p \times 1}$  with  $h_j := (\kappa_k'(w_j))_{k=1}^K$

target function:  $f(x) = \sum_{k=1}^K \alpha_k \kappa_k(x)$

$$R^{\text{RFM}} := E_{\mathcal{D}} \left[ \frac{1}{p} H^T G^{-2} \right]$$

## Background: Test risk of random features model

A compact integral operator  $T : L^2(\mathcal{X}) \rightarrow L^2(\mathcal{W})$  for any  $f \in L^2(\mathcal{X})$   
(De Lippis et al., 2024)

$$(Tf)(w) := \int_{\mathcal{R}^d} \kappa(x; w) f(x) d(x); \quad T = \sum_{k=1}^K \kappa_k \kappa_k'$$

Covariate feature matrix  $G := [g_1; \dots; g_n] \in \mathbb{R}^{n \times 1}$  with  $g_i := (\kappa_k(x_i))_{k=1}^K$

Weight feature matrix  $H := [h_1; \dots; h_p] \in \mathbb{R}^{p \times 1}$  with  $h_j := (\kappa_k'(w_j))_{k=1}^K$

target function:  $f(x) = \sum_{k=1}^K \alpha_k \kappa_k(x)$

$$R^{\text{RFM}} := E_{\mathcal{D}} \left[ \frac{1}{p} H^T \hat{a} \right]^2$$

A compact integral operator  $T : L^2(\mathcal{X}) \rightarrow L^2(\mathcal{W})$  for any  $f \in L^2(\mathcal{X})$   
 (De Lippis et al., 2024)

$$(Tf)(w) := \int_{\mathcal{R}^d} \kappa(x; w) f(x) d(x); \quad T = \sum_{k=1}^K \kappa_k \kappa_k'$$

Covariate feature matrix  $G := [g_1; \dots; g_n] \in \mathbb{R}^{n \times 1}$  with  $g_i := (\kappa_k(x_i))_{k=1}^K$

Weight feature matrix  $H := [h_1; \dots; h_p] \in \mathbb{R}^{p \times 1}$  with  $h_j := (\kappa_k'(w_j))_{k=1}^K$

target function:  $f(x) = \sum_{k=1}^K \alpha_k \kappa_k(x)$

$$R^{RFM} := E \left[ \frac{1}{p} H^T G^{-2} \right]$$

## Background: Test risk of random features model

A compact integral operator  $T : L^2(\mathcal{X}) \rightarrow L^2(\mathcal{W})$  for any  $f \in L^2(\mathcal{X})$   
(De Lippis et al., 2024)

$$(Tf)(w) := \int_{\mathcal{R}^d} \kappa(x; w) f(x) d(x); \quad T = \sum_{k=1}^K \kappa_k \kappa_k'$$

Covariate feature matrix  $G := [g_1; \dots; g_n] \in \mathbb{R}^{n \times 1}$  with  $g_i := (\kappa_k(x_i))_{k=1}^K$

Weight feature matrix  $H := [h_1; \dots; h_p] \in \mathbb{R}^{p \times 1}$  with  $h_j := (\kappa_k'(w_j))_{k=1}^K$

target function:  $f(x) = \sum_{k=1}^K \alpha_k \kappa_k(x)$

$$R^{\text{RFM}} := E_n \left[ \frac{1}{p} H^T a \right]^2$$

# Empirical observation under real-world dataset

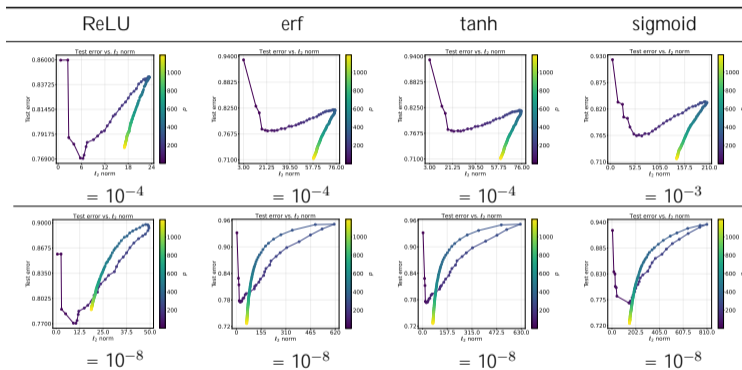


Figure 3: Results for RFMs under FasionMNIST.

# Our results under well-behaved data

(a) Test Risk vs.  $\ell_2$  norm vs.  $\ell_1$  norm (b)  $\ell_2$  norm vs.  $\ell_1$  norm (c) Test Risk vs. norm (d)  $\epsilon = 0.001$

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

# Our results under well-behaved data

(a) Test Risk vs.  $\ell_2$  norm vs.  $\ell_1$  norm (b)  $\ell_2$  norm vs.  $\ell_1$  norm (c) Test Risk vs.  $\ell_2$  norm (d)  $\rho = 0.001$

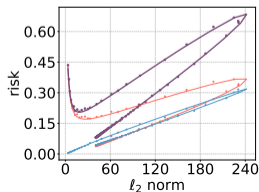
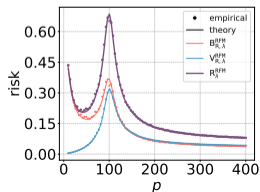
- $\rho := p/n$ ,  $p$ : model size (width),  $n$ : data size
- Phase transition exists but double descent does not exist
- More close to U-shaped instead of double descent: **A ' paradigm**
- Over-parameterization is still better than under-parameterization

# Our results under well-behaved data

(a) Test Risk vs.  $p$       (b)  $\ell_2$  norm vs.  $p$       (c) Test Risk vs.  $\ell_2$  norm      (d)  $\sigma^2 = 0.001$

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

Test error = Bias<sup>2</sup> + Variance



# Our results under well-behaved data

(a) Test Risk vs.  $\ell_2$  norm vs.  $\ell_2$  norm (d) = 0.001

- $p := p=n$ ,  $p$ : model size (width),  $n$ : data size
- Reshape scaling-law:  
test loss =  $A \text{ Data Size}^a + B \text{ Model Size}^b + C$  with  $a; b > 0$

# Our results under well-behaved data

(a) Test Risk vs.  $\ell_2$  norm vs.  $\ell_2$  norm (d) = 0.001

- $p = n$ ,  $p$ : model size (width),  $n$ : data size
- Reshape scaling-law:  
test loss = A Data Size<sup>a</sup> + B Model Size<sup>b</sup> + C with  $a, b > 0$   
test loss = A Data Size<sup>a</sup> Norm Capacity<sup>b</sup> with  $a > 0$  and  $b \geq R$

## Control norm by tuning $\lambda$ : L-curve (Hansen, 1992)

Explicit (model size) vs. Implicit (norm)

One-to-one mapping between norm and

## Control norm by tuning $\lambda$ : L-curve (Hansen, 1992)

Explicit (model size) vs. Implicit (norm)

One-to-one mapping between norm and

(a) Norm vs. (varying )      (b) Risk vs. Norm (varying )

## An example of linear regression: Textbook level and beyond

- $n$  i.i.d. samples  $f(x_i; y_i)_{i=1}^n$  with  $x_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$
- $y = h(x) + \epsilon$ ,  $E(\epsilon) = 0$  and  $V(\epsilon) = \sigma^2$ , covariance matrix  $\Sigma = E[xx^T]$
- ridge regression:  $\hat{\beta} = (X^T X + \lambda I)^{-1} X^T y$

Target: precise analysis

The expected test risk  $E \|\hat{f} - f\|_{L_2}^2 \sim k^2$  vs. the norm  $E \|\hat{f} - k\|_{L_2}^2$

## An example of linear regression: Textbook level and beyond

- $n$  i.i.d. samples  $f(x_i; y_i)_{i=1}^n$  with  $x_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$
- $y = h(x) + \epsilon$ ,  $E(\epsilon) = 0$  and  $V(\epsilon) = \sigma^2$ , covariance matrix  $\Sigma = E[xx^T]$
- ridge regression:  $\hat{\beta} = (X^T X + \lambda I)^{-1} X^T y$

Target: precise analysis

The expected test risk  $E \|\hat{h} - h\|_{k^2}^2$  vs. the norm  $E \|\hat{h}\|_{k^2}^2$

# An example of linear regression: Textbook level and beyond

- $n$  i.i.d. samples  $f(x_i; y_i) g_{i=1}^n$  with  $x_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$
- $y = h(x) + \epsilon$ ,  $E(\epsilon) = 0$  and  $V(\epsilon) = \sigma^2$ , covariance matrix  $\Sigma = E[xx^T]$
- ridge regression:  $\hat{\beta} = (X^T X + \lambda I)^{-1} X^T y$

Target: precise analysis

The expected test risk  $E \|\hat{k} - k\|^2$  vs. the norm  $E \|\hat{k}\|_2^2$

- Deterministic equivalence (Cheng and Montanari, 2024; Misiakiewicz and Saeed, 2024; Bach, 2024)

The empirical spectral measure converges to a deterministic limit.

# An example of linear regression: Textbook level and beyond

- $n$  i.i.d. samples  $f(x_i; y_i) g_{i=1}^n$  with  $x_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$
- $y = h(x) + \epsilon$ ,  $E(\epsilon) = 0$  and  $V(\epsilon) = \sigma^2$ , covariance matrix  $\Sigma = E[xx^T]$
- ridge regression:  $\hat{\beta} = (X^T X + \lambda I)^{-1} X^T y$

Target: precise analysis

The expected test risk  $E \|\hat{\beta} - \beta\|^2$  vs. the norm  $E \|\hat{\beta}\|_2^2$

- Deterministic equivalence (Cheng and Montanari, 2024; Misiakiewicz and Saeed, 2024; Bach, 2024)

$$\text{Tr} X^T X (X^T X + \lambda I)^{-1} = \text{Tr} (\Sigma + \lambda I)^{-1}; w:h:p:$$

- can be asymptotic or non-asymptotic at the rate of  $O(1 - \frac{p}{n})$ .
- $\lambda$  is the non-negative solution to the self-consistent equation  $n - p = \text{Tr}(\Sigma + \lambda I)^{-1}$ .

# An example of linear regression: Textbook level and beyond

- $n$  i.i.d. samples  $f(x_i; y_i) g_{i=1}^n$  with  $x_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$
- $y = h(x) + \epsilon$ ,  $E(\epsilon) = 0$  and  $V(\epsilon) = \sigma^2$ , covariance matrix  $\Sigma = E[xx^T]$
- ridge regression:  $\hat{\beta} = (X^T X + \lambda I)^{-1} X^T y$

Target: precise analysis

The expected test risk  $E \|\hat{k} - k\|^2$  vs. the norm  $E \|\hat{k}\|_2^2$

- Deterministic equivalence (Cheng and Montanari, 2024; Misiakiewicz and Saeed, 2024; Bach, 2024)
- Bias-variance decomposition on the test risk
  - $B_{R;LS}^2 = \sigma^2 h^T (X^T X + \lambda I)^{-1} (X^T X + \lambda I)^{-1} i$
  - $V_{R;LS}^2 = \sigma^2 \text{Tr}(X^T X (X^T X + \lambda I)^{-2})$

## Theorem (asymptotic/non-asymptotic results)

We have a bias-variance decomposition  $E \|k^{\wedge} k_2^2 = B_{N_n} + V_{N_n}$ .

For *well-behaved* data and  $\lambda$ , we have  $B_{N_n} \rightarrow B_{N_n}$  and  $V_{N_n} \rightarrow V_{N_n}$ ; w.h.p.

$$B_{N_n} := \frac{2h}{n} \text{Tr} \left( \frac{(\lambda + I)^{-2}}{n} \right) + \frac{\text{Tr} \left( (\lambda + I)^{-2} \right)}{n} \frac{2h}{1 - \frac{1}{n} \text{Tr} \left( \frac{(\lambda + I)^{-2}}{n} \right)};$$

|-----{Z}-----|  
 $B_{N_n}$

$$V_{N_n} := \frac{2 \text{Tr} \left( (\lambda + I)^{-2} \right)}{n \text{Tr} \left( \frac{(\lambda + I)^{-2}}{n} \right)};$$

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

- Norm-based capacity:  $3 \lambda$ ,
- effective dimension-style  $\text{Tr} \left( (\lambda + I)^{-1} \right)$ :  $7 \lambda$





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

• Test risk  $R$  and norm  $N$  formulates a cubic curve (complex but precise).

• min-norm interpolator ( $\lambda = 0$ ):

$$R_0 = \frac{N_0}{N_0 + (k - k_2^2)^2 + 4k k_2^2} \quad ; \quad \text{in under-parameterized regimes}$$

• optimal regularization  $\lambda = \frac{d^2}{k - k_2^2}$  (Wu and Xu,

$$2020): R = k - k_2^2 - N$$

•  $\lambda \rightarrow \infty$ :  $R = k - k_2^2 - \frac{d^2}{N}$

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

• Test risk  $R$  and norm  $N$  formulates a cubic curve (complex but precise).

- min-norm interpolator ( $\lambda = 0$ ):

$$R_0 = \frac{N_0}{N_0 + (k - k_2^2)^2 + 4k k_2^2} \quad ; \quad \text{in under-parameterized regimes}$$

- optimal regularization  $\lambda = \frac{d^2}{k - k_2^2}$  (Wu and Xu, 2020):  $R = k - k_2^2 - N$

- $\lambda \rightarrow \infty$ :  $R = k - k_2^2 - \frac{\rho}{N^2}$

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

• Test risk  $R$  and norm  $N$  formulates a cubic curve (complex but precise).

- min-norm interpolator ( $\lambda = 0$ ):

$$R_0 = \frac{N_0}{N_0 + (k - k_2^2)^2 + 4k k_2^2} \quad ; \quad N_0 = k k_2^2$$

Why? Variance is the same

= 0 (under-parameterized)

=  $\frac{d}{n}$  (over-parameterized)

- optimal regularization  $\lambda = \frac{d}{k - k_2^2}$  (Wu and Xu, 2020):  $R = k - k_2^2 - N$

- $\lambda \rightarrow 1$ :  $R = k - k_2^2 - \frac{d}{N}$

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

• Test risk  $R$  and norm  $N$  formulates a cubic curve (complex but precise).

• min-norm interpolator ( $\lambda = 0$ ):

$$R_0 = \frac{8}{9} \frac{N_0 k k_2^2}{N_0 (k k_2^2)^2 + 4k k_2^2} : \quad \text{in under-parameterized regimes}$$

Why? Variance is the same

= 0 (under-parameterized)

=  $\frac{d}{n}$  (over-parameterized)

• optimal regularization  $\lambda = \frac{d}{k k_2^2}$  (Wu and Xu,

2020):  $R = k k_2^2 N$

•  $\lambda = 1$ :  $R = k k_2^2 \frac{d}{N^2}$

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

$\rho$  Test risk  $R$  and norm  $N$  formulates a cubic curve (complex but precise).

- min-norm interpolator ( $\rho = 0$ ):

$$R_0 = \frac{8}{9} \frac{N_0 k k_2^2}{N_0 (k k_2^2)^2 + 4k k_2^2} : \quad \text{in under-parameterized regimes}$$

Why? Variance is the same

$= 0$  (under-parameterized)

$= \frac{d}{n}$  (over-parameterized)

- optimal regularization  $\rho = \frac{d}{k k_2^2}$  (Wu and Xu,

2020):  $R = k k_2^2 N$

- $\rho = 1$ :  $R = k k_2^2 \frac{\rho}{N}^2$

## 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.

## 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.

Is  $\ell_2$  norm-based capacity best for characterizing generalization?

# 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	$\sum_{i=1}^L \ W_i\ _2$	0:537
Fisher-Rao	$\frac{(L+1)^2}{n} \sum_{i=1}^n \sum_{j=1}^L \langle h_{W_j}; \nabla_{W_j} \langle h_W(x_i); y_i \rangle \rangle$	0:078
Path-norm	$\sum_{j=1}^L \ W_j\ _1$	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	$\sum_{i=1}^L \ W_i\ _F^2$	0:073
Frobenius distance to initialization	$\sum_{i=1}^L \ W_i - W_i^0\ _F^2$	0:263
Spectral complexity	$\sum_{i=1}^L \ W_i\ _2^3$	0:537
Fisher-Rao	$\frac{(L+1)^2}{n} \sum_{i=1}^n \sum_{j=1}^L \langle \nabla_{W_j} \ell(x_i; y_i), \nabla_{W_j} \ell(x_i; y_i) \rangle$	0:078
Path-norm	$\sum_{j=1}^L \sum_{i=1}^n \ W_j(x_i)\ _2$	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	$\sum_{i=1}^L \ W_i\ _F^2$	0:073
Frobenius distance to initialization	$\sum_{i=1}^L \ W_i - W_i^0\ _F^2$	0:263
Spectral complexity	$\sum_{i=1}^L \ W_i\ _k^{\frac{3=2}{k=2}}$	0:537
Fisher-Rao	$\frac{(L+1)^2}{n} \sum_{i=1}^n \langle h_{W_i}; r_{W_i} \rangle$	0:078
Path-norm	$\sum_{(i_0, \dots, i_L)} \sum_{j=1}^L \ W_{i_j; i_{j-1}}\ _2$	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	$\sum_{i=1}^L \ W_i\ _F^2$	0:073
Frobenius distance to initialization	$\sum_{i=1}^L \ W_i - W_i^0\ _F^2$	0:263
Spectral complexity	$\sum_{i=1}^L \ W_i\ _2^3$	0:537
Fisher-Rao	$\frac{(L+1)^2}{n} \sum_{i=1}^n \sum_{j=1}^L \langle \nabla_{W_j} \ell(x_i; y_i), \nabla_{W_j} \ell(x_i; y_i) \rangle$	0:078
Path-norm	$\sum_{j=1}^L \sum_{i_0, \dots, i_L} \ W_{i_j}\ _2$	0:373

(a) Test (training) Loss vs.  $p$       (b) Path-norm vs.  $p$       (c) Test Loss vs. Path-norm

Figure 6: Experiments on two-layer neural networks.

# Which model capacity is suitable (for neural networks)?

Table 2: 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	$\sum_{i=1}^L \ W_i\ _2$	0:537
Fisher-Rao	$\frac{(L+1)^2}{n} \sum_{i=1}^n \sum_{j=1}^L \langle h_{W_j}; \nabla_{W_j} \langle h_{W_j}(x_i); y_i \rangle \rangle$	0:078
Path-norm	$\sum_{j=1}^L \ W_j\ _1$	0:373

# Which model capacity is suitable (for neural networks)?

Table 2: 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	$\sum_{i=1}^L \ W_i\ _2^2$	0:537
Fisher-Rao	$\frac{(L+1)^2}{n} \sum_{i=1}^n \sum_{j=1}^L \langle \nabla_{W_j} h_W(x_i); \nabla_{W_j} h_W(x_i) \rangle$	0:078
Path-norm	$\sum_{j=1}^L \sum_{i=1}^n \ W_j \cdot \nabla_{W_j} h_W(x_i)\ _2$	0:373

# Which model capacity is suitable (for neural networks)?

Table 2: 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	$\sum_{i=1}^L \ W_i\ _k^2$ $\sum_{i=1}^L \frac{\ W_i\ _{2,1}^{3=2}}{\ W_i\ _{3=2}^{3=2}}$ $! \quad 2=3$	0:537
Fisher-Rao	$\frac{(L+1)^2}{n} \sum_{i=1}^n \langle h_{W_i}; r_{W_i} \rangle \langle h_{W_i}(x_i); y_i \rangle$	0:078
Path-norm	$\sum_{(i_0, \dots, i_L)} \sum_{j=1}^L \ W_{i_j; i_{j-1}}\ _1^2$	0:373

# Which model capacity is suitable (for neural networks)?

Table 2: 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	$\sum_{i=1}^L \ W_i\ _2^3$	0:537
Fisher-Rao	$\frac{(L+1)^2}{n} \sum_{i=1}^n \sum_{j=1}^L \langle \nabla_{W_j} \ell, \nabla_{W_j} \ell \rangle$	0:078
Path-norm	$\sum_{j=1}^L \sum_{i_0, \dots, i_L} \ W_{i_j}\ _2$	0:373

(a) Test (training) Loss vs.  $p$       (b) Path-norm vs.  $p$       (c) Test Loss vs. Path-norm

Figure 7: Experiments on two-layer neural networks.

## Two-layer neural networks, path norm

---

$\ell_1$ -path norm (Neyshabur et al., 2015)

$$\|k\|_{k_P} := \frac{1}{m} \sum_{k=1}^m |a_k| w_k$$

$\ell_1$ -path norm (Neyshabur et al., 2015)

$$\|k\|_{k_P} := \frac{1}{m} \sum_{k=1}^m \|a_k\|_{W_k} \|k_1\|$$

- equivalent to Barron spaces  $B$  (Barron, 1993; E et al., 2021)

$$B := \{ f_a : \|a\|_{L^2(\cdot)} < 1 \}$$

$\ell_1$ -path norm (Neyshabur et al., 2015)

$$\|k\|_{k_P} := \frac{1}{m} \sum_{k=1}^m \|a_k\|_{W_k} \|k\|_1$$

- equivalent to Barron spaces  $B$  (Barron, 1993; E et al., 2021)

$$B := \{ f_a : \|a\|_{L^2(\cdot)} < 1 \}$$

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

$\ell_1$ -path norm (Neyshabur et al., 2015)

$$\|k\|_{\ell_1} := \frac{1}{m} \sum_{k=1}^m \|a_k\|_{L^2(\cdot)} \|W_k\|_{L^2(\cdot)}$$

- equivalent to Barron spaces  $B$  (Barron, 1993; E et al., 2021)

$$B := \{ f : \mathbb{R}^d \rightarrow \mathbb{R} : \|f\|_B < 1 \}$$

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

*Can neural networks identify this structure?*

Theorem (Informal, sample complexity of learning  $f^? \geq B$ )

To achieve  $\epsilon$ -excess risk,

- Kernel methods require  $(\frac{1}{\epsilon})^d$  samples.
- Two-layer neural networks require  $(\frac{2d+2}{\epsilon})$  samples. *smaller than  $(\frac{1}{\epsilon})^2$*

To achieve  $\epsilon$ -excess risk,

^ Kernel methods require  $\binom{d}{2}$  samples.

^ Two-layer neural networks require  $\binom{2d+2}{d+2}$  samples, smaller than  $d^2$

To achieve  $\epsilon$ -excess risk,

^ Kernel methods require  $(\epsilon^{-d})$  samples.

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

No Curse of Dimensionality: NNs adapt to directional smoothness.

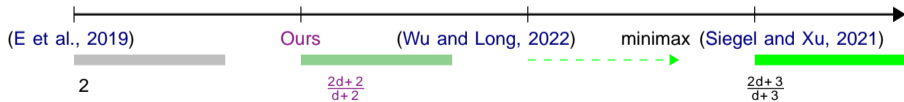
To achieve  $\epsilon$ -excess risk,

^ Kernel methods require  $(\epsilon^{-d})$  samples.

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

No Curse of Dimensionality: NNs adapt to directional smoothness.

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



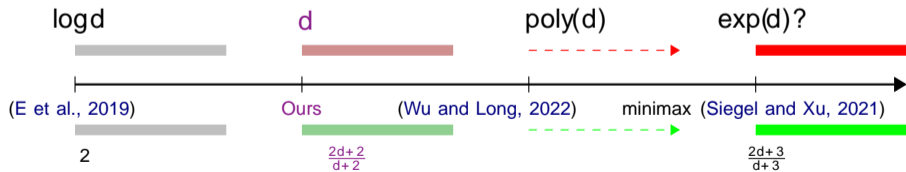
To achieve  $\epsilon$ -excess risk,

^ Kernel methods require  $(d)$  samples.

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

No Curse of Dimensionality: NNs adapt to directional smoothness.

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



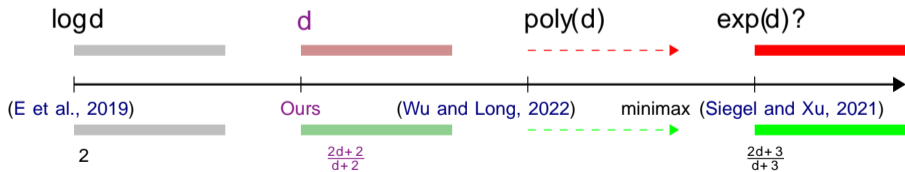
To achieve  $\epsilon$ -excess risk,

^ Kernel methods require  $(d)$  samples.

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

No Curse of Dimensionality: NNs adapt to directional smoothness.

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



The best trade-o between  $\epsilon$  and  $d$ .

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 ) interesting mathematical problems

p Be aware of model capacity! **A new paradigm'ofcurve!**

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

Deep learning phenomena ) interesting mathematical problems

- p Be aware of model capacity! **A new paradigm of curve!**
  - ^ Reshape bias-variance trade-offs, double descent, scaling law under proper norm-based capacity via **deterministic equivalence**.
  
- p Which function class can be efficiently learned by neural networks?
  - ^ Neural networks can adapt to low-dimensional structure and avoid CoD!

### Deep learning phenomena ) interesting mathematical problems

- p Be aware of model capacity! **A new paradigm of curve!**
  - ^ Reshape bias-variance trade-offs, double descent, scaling law under proper norm-based capacity via **deterministic equivalence**.

- p Which function class can be efficiently learned by neural networks?
  - ^ Neural networks can adapt to low-dimensional structure and avoid CoD!

### Theoretical advances ) principled guidance in practical problems

- p How does theory contribute to practical fine-tuning problems?
  - ^ One-step full gradient can be sufficient! **[ICML'25 oral]**

## 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.
- Francis Bach. High-dimensional analysis of double descent for linear regression with random projections. *SIAM Journal on Mathematics of Data Science*, 6(1):26–50, 2024.

- 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.
- Hongrui Chen, Jihao Long, and Lei Wu. A duality framework for generalization analysis of random feature models and two-layer neural networks. *arXiv preprint arXiv:2305.05642* 2023.

- 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. *Annals of Statistics* 52(6):2879–2912, 2024.
- Leonardo De Lippis, Bruno Loureiro, and Theodor Misiakiewicz. Dimension-free deterministic equivalents for random feature regression. In *Advances in Neural Information Processing Systems*, 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 low-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. *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.

- Fanghui Liu, Xiaolin Huang, Yudong Chen, and Johan AK Suykens. Random features for kernel approximation: A survey on algorithms, theory, and beyond. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44(10):7128–7148, 2021.
- 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 tradeo : 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](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 finite 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  $\epsilon$ -widths of shallow neural networks. *arXiv preprint arXiv:2101.12365*, 2021.
- Aad W Van Der Vaart, Adrianus Willem van der Vaart, Aad van der Vaart, and Jon Wellner. *Weak convergence and empirical processes: with applications to statistics* Springer Science & Business Media, 1996.

- 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 regularization in overparameterized linear regression. 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.

(a) Test (training) Loss vs.  $p$       (b) Fro-norm vs.  $p$       (c) Test Loss vs. Fro-norm

Figure 8: Experiments on two-layer fully connected neural networks with noise level  $\sigma = 0.2$ . The left figure shows the relationship between test (training) loss and the number of the parameters  $p$ . The middle figure shows the relationship between the Frobenius norm and  $p$ . The right figure shows the relationship between the test loss and Fro-norm.

## An example of linear model: a textbook level

- $f(x_i; y_i) g_{i=1}^n$  i.i.d.: ,  $x_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$ , covariance matrix  $\Sigma = E[xx^T]$
- $y = h^T x + \epsilon$  with  $E(\epsilon) = 0$  and  $V(\epsilon) = \sigma^2$
- ridge regression:  $\hat{\beta} = (X^T X + \lambda I)^{-1} X^T y$
- min- $\ell_2$ -norm interpolation:  $\hat{\beta}_{\min} = \operatorname{argmin}_{\beta} \|\beta\|_2; \text{s.t. } X\beta = y$
- expected test risk: bias-variance decomposition

$$R^{\text{LS}} := E \|\hat{\beta} - \beta^*\|^2 = \underbrace{E \|\hat{\beta} - \beta^*\|^2}_{:= B_R^{\text{LS}}} + \underbrace{\operatorname{tr}(\operatorname{Cov}(\hat{\beta}))}_{:= V_R^{\text{LS}}}$$

- $B_R^{\text{LS}} = \sigma^2 h^T (X^T X + \lambda I)^{-1} (X^T X + \lambda I)^{-1} h$
- $V_R^{\text{LS}} = \sigma^2 \operatorname{Tr}(X^T X (X^T X + \lambda I)^{-2})$
- \*Intuitive fact: for i.i.d. sub-Gaussian data  $X$ , we have

$$\frac{1}{n} X^T X \approx \frac{p}{d-n} I; w.h.p.:$$

## An example of linear model: a textbook level

- $f(x_i; y_i) g_{i=1}^n$  i.i.d.,  $x_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$ , covariance matrix  $\Sigma = E[xx^T]$
- $y = h^T x + \epsilon$  with  $E(\epsilon) = 0$  and  $V(\epsilon) = \sigma^2$
- ridge regression:  $\hat{\beta} = (X^T X + \lambda I)^{-1} X^T y$
- min- $\ell_2$ -norm interpolation:  $\hat{\beta}_{\min} = \arg \min_k \|k\|_2; \text{s.t. } X \hat{\beta} = y$
- expected test risk: bias-variance decomposition

$$R^{\text{LS}} := E \|k\|^2 \quad \hat{k}^2 = \underbrace{\|E[\hat{\beta}]k\|^2}_{:= B_R^{\text{LS}}} + \underbrace{\text{tr}(\text{Cov}(\hat{\beta}))}_{:= V_R^{\text{LS}}}$$

- $B_R^{\text{LS}} = \sigma^2 h^T (X^T X + \lambda I)^{-1} (X^T X + \lambda I)^{-1} h$
- $V_R^{\text{LS}} = \sigma^2 \text{Tr}(X^T X (X^T X + \lambda I)^{-2})$
- \*Intuitive fact: for i.i.d. sub-Gaussian data  $X$ , we have

$$\frac{1}{n} X^T X \approx \frac{p}{d-n} I; w.h.p.$$

## An example of linear model: a textbook level

- $f(x_i; y_i) g_{i=1}^n$  i.i.d.,  $x_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$ , covariance matrix  $\Sigma = E[xx^T]$
- $y = h^T x + \epsilon$  with  $E(\epsilon) = 0$  and  $V(\epsilon) = \sigma^2$
- ridge regression:  $\hat{\beta} = (X^T X + \lambda I)^{-1} X^T y$
- min- $\ell_2$ -norm interpolation:  $\hat{\beta}_{\min} = \arg \min_k \|k\|_2; \text{s.t. } X \hat{\beta} = y$
- expected test risk: bias-variance decomposition

$$R^{\text{LS}} := E \|k\|^2 = \underbrace{\|E[\hat{\beta}]\|_2^2}_{:= B_R^{\text{LS}}} + \underbrace{\text{tr}(\text{Cov}(\hat{\beta}))}_{:= V_R^{\text{LS}}}$$

- $B_R^{\text{LS}} = \sigma^2 h^T (X^T X + \lambda I)^{-1} (X^T X + \lambda I)^{-1} h$
- $V_R^{\text{LS}} = \sigma^2 \text{Tr}(X^T X (X^T X + \lambda I)^{-2})$
- \*Intuitive fact: for i.i.d. sub-Gaussian data  $X$ , we have

$$\frac{1}{n} X^T X \approx \frac{p}{d-n} I; w.h.p.$$

## An example of linear model: a textbook level

- $f(x_i; y_i) g_{i=1}^n$  i.i.d.,  $x_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$ , covariance matrix  $\Sigma = E[xx^T]$
- $y = h^T x + \epsilon$  with  $E(\epsilon) = 0$  and  $V(\epsilon) = \sigma^2$
- ridge regression:  $\hat{\beta} = (X^T X + \lambda I)^{-1} X^T y$
- min- $\ell_2$ -norm interpolation:  $\hat{\beta}_{\min} = \arg \min_k \|k\|_2; \text{s.t. } X \hat{\beta} = y$
- expected test risk: bias-variance decomposition

$$R^{\text{LS}} := E \|k\|^2 = \underbrace{\|E[\hat{\beta}]\|_2^2}_{:= B_R^{\text{LS}}} + \underbrace{\text{tr}(\text{Cov}(\hat{\beta}))}_{:= V_R^{\text{LS}}}$$

- $B_R^{\text{LS}} = \sigma^2 h^T (X^T X + \lambda I)^{-1} (X^T X + \lambda I)^{-1} h$
- $V_R^{\text{LS}} = \sigma^2 \text{Tr}(X^T X (X^T X + \lambda I)^{-2})$
- \*Intuitive fact: for i.i.d. sub-Gaussian data  $X$ , we have

$$\frac{1}{n} X^T X \approx \frac{p}{d-n} I; w.h.p.:$$

# Beyond textbook level: deterministic equivalence (Cheng and Montanari, 2024)

$$\text{Tr} X^>X(X^>X + I)^{-1} = \text{Tr} (I + I)^{-1} :$$

- can be asymptotic or non-asymptotic at the rate of  $O(1/\sqrt{n})$ .
- is the non-negative solution to the self-consistent equation  $n^{-1} \text{Tr} (I + I)^{-1}$ .

Theorem (Deterministic equivalence (Cheng and Montanari, 2024))

For sub-Gaussian data, assume  $\Sigma$  is well-behaved, w.h.p.

$$\frac{\frac{1}{n} \text{Tr} \{ \hat{Z}^k \}}{k} := B_R^{\text{LS}} \quad B_R^{\text{LS}} := \frac{1}{n} \text{tr} ( \Sigma (I + \Sigma)^{-2} )$$

$$\frac{\text{tr}(\text{Cov}(\hat{Z}))}{n} := V_R^{\text{LS}} \quad V_R^{\text{LS}} := \frac{2 \text{tr} ( \Sigma (I + \Sigma)^{-2} )}{n \text{tr} ( \Sigma (I + \Sigma)^{-2} )}$$

# Beyond textbook level: deterministic equivalence (Cheng and Montanari, 2024)

$$\text{Tr} X^>X(X^>X + I)^{-1} = \text{Tr} (X^>X + I)^{-1} :$$

- can be asymptotic or non-asymptotic at the rate of  $O(1/\sqrt{n})$ .
- is the non-negative solution to the self-consistent equation  $n^{-1} \text{Tr} (X^>X + I)^{-1}$ .

Theorem (Deterministic equivalence (Misiakiewicz and Saeed, 2024))

For sub-Gaussian data, assume  $X$  is well-behaved, w.h.p.

$$\underbrace{\frac{1}{n} \text{Tr} \{Z^>Z\}}_{:= B_R^{LS}} := \frac{\sigma^2}{1 + n^{-1} \text{tr}(\sigma^2 (X^>X + I)^{-2})}$$

$$\underbrace{\frac{1}{n} \text{Tr} \{ \text{Cov}(Z^>Z) \}}_{:= V_R^{LS}} := \frac{\sigma^2 \text{tr}(\sigma^2 (X^>X + I)^{-2})}{n \text{tr}(\sigma^2 (X^>X + I)^{-2})} :$$

## \*Path norm, Barron spaces, RKHS (Chen et al., 2023)

Consider a random features model (RFM) (Rahimi and Recht, 2007)

- first layer:  $w \stackrel{iid}{\sim} P(W)$ ; only train the second layer

infinite many features  $f_a(x) = \int_W a(w) \phi(x; w) d(w)$

$$F_p := \{f_a : \|a\|_{L^p(\cdot)} < 1\}; \quad \|f\|_{F_p} := \inf_{f=f_a} \|a\|_{L^p(\cdot)}$$

- RFMs  $\approx$  kernel methods by taking  $p = 2$  using Representer theorem
- RFMs  $\not\approx$  kernel methods if  $p < 2$
- function space:  $F_1 \subset F_p \subset F_q \subset F_1$ ; if  $p < q$

For any  $1 \leq p < \infty$ , define

$$B = \left[ \int_{2P(W)} F_p; \quad \|f\|_B = \inf_{f=f_a} \|a\|_{L^p(\cdot)} \right]$$

largest

data-adaptive

## \*Path norm, Barron spaces, RKHS (Chen et al., 2023)

Consider a random features model (RFM) (Rahimi and Recht, 2007)

- first layer:  $w \stackrel{iid}{\sim} P(W)$ ; only train the second layer

infinite many features  $f_a(x) = \int_{\mathcal{W}} a(w) \phi(x; w) d(w)$

$$F_p := \{f_a : \|a\|_{L^p(\cdot)} < 1\}; \quad \|f\|_{F_p} := \inf_{f=f_a} \|a\|_{L^p(\cdot)}$$

- RFMs  $\approx$  kernel methods by taking  $p = 2$  using Representer theorem
- RFMs  $\neq$  kernel methods if  $p < 2$
- function space:  $F_1 \subset F_p \subset F_q \subset F_1$ ; if  $p < q$

For any  $1 \leq p < \infty$ , define

$$B = \left[ \int_{\mathcal{W}} P(W) F_p; \quad \|f\|_B = \inf_{f=f_a} \|a\|_{F_p} \right]$$

largest  
data-adaptive

## \*Path norm, Barron spaces, RKHS (Chen et al., 2023)

Consider a random features model (RFM) (Rahimi and Recht, 2007)

- first layer:  $w \stackrel{iid}{\sim} P(W)$ ; only train the second layer

infinite many features  $f_a(x) = \int_{\mathcal{W}} a(w) \phi(x; w) d(w)$

$$F_{p; \phi} := \{f_a : \|a\|_{L^p(\mathcal{W})} < 1\}; \quad \|f\|_{F_{p; \phi}} := \inf_{f=f_a} \|a\|_{L^p(\mathcal{W})}$$

- RFMs  $\approx$  kernel methods by taking  $p = 2$  using Representer theorem
- RFMs  $\not\approx$  kernel methods if  $p < 2$
- function space:  $F_1 \subset F_p \subset F_q \subset F_1$ ; if  $p < q$

For any  $1 \leq p < \infty$ , define

$$B = \left[ \int_{\mathcal{W}} \phi(x; w)^2 d(w) \right]^{-1/2} F_p; \quad \|f\|_B = \inf_{f=f_a} \|a\|_{L^p(\mathcal{W})}$$

largest  
data-adaptive

## \*Path norm, Barron spaces, RKHS (Chen et al., 2023)

Consider a random features model (RFM) (Rahimi and Recht, 2007)

- first layer:  $w \stackrel{iid}{\sim} P(W)$ ; only train the second layer

infinite many features  $f_a(x) = \int_{\mathcal{W}} a(w) \phi(x; w) d(w)$

$$F_{p; \gamma} := \{f_a : \|a\|_{L^p(\gamma)} < \gamma\}; \quad \|f\|_{F_{p; \gamma}} := \inf_{f=f_a} \|a\|_{L^p(\gamma)}$$

- RFMs  $\approx$  kernel methods by taking  $p = 2$  using Representer theorem
- RFMs  $\not\approx$  kernel methods if  $p < 2$
- function space:  $F_1 \subset F_p \subset F_q \subset F_1$ ; if  $p < q$

For any  $1 \leq p < \infty$ , define

$$B = \{f \in L^2(\gamma) : \|f\|_B = \inf_{f=f_a} \|a\|_{L^p(\gamma)}\}$$

largest  
data-adaptive

## \*Path norm, Barron spaces, RKHS (Chen et al., 2023)

Consider a random features model (RFM) (Rahimi and Recht, 2007)

- first layer:  $w \stackrel{iid}{\sim} P(W)$ ; only train the second layer

infinite many features  $f_a(x) = \int_{\mathcal{W}} a(w) \phi(x; w) d(w)$

$$F_{p; \gamma} := \{f_a : \|a\|_{L^p(\gamma)} < \gamma\}; \quad \|f\|_{F_{p; \gamma}} := \inf_{f=f_a} \|a\|_{L^p(\gamma)}$$

- RFMs  $\approx$  kernel methods by taking  $p = 2$  using Representer theorem
- RFMs  $\not\approx$  kernel methods if  $p < 2$
- function space:  $F_1 \subset F_p \subset F_q \subset F_1$ ; if  $p < q$

For any  $1 \leq p < \infty$ , define

$$B = \left[ \int_{\mathcal{W}} P(W) F_{p; \gamma} \right]; \quad \|f\|_B = \inf_{f=f_a} \|a\|_{F_p}$$

largest  
data-adaptive

## \*Path norm, Barron spaces, RKHS (Chen et al., 2023)

Consider a random features model (RFM) (Rahimi and Recht, 2007)

- first layer:  $w \stackrel{iid}{\sim} P(W)$ ; only train the second layer

infinite many features  $f_a(x) = \int_{\mathcal{W}} a(w) \phi(x; w) d(w)$

$$F_{p; \gamma} := \{f_a : \|a\|_{L^p(\gamma)} < \gamma\}; \quad \|f\|_{F_{p; \gamma}} := \inf_{f=f_a} \|a\|_{L^p(\gamma)}$$

- RFMs  $\approx$  kernel methods by taking  $p = 2$  using Representer theorem
- RFMs  $\not\approx$  kernel methods if  $p < 2$
- function space:  $F_1 \subset F_p \subset F_q \subset F_1$ ; if  $p < q$

For any  $1 \leq p < \infty$ , define

$$B = \left[ \int_{\mathcal{W}} \phi(x; w)^2 d(w) \right]^{-1/2} F_p; \quad \|f\|_B = \inf_{f=f_a} \|a\|_{F_p}$$

largest

data-adaptive

# Proof sketch: convex hull technique and its constant!

- Consider the following function space

$$F = \{f : \mathbb{R}^d \rightarrow \mathbb{R} \mid f(x) = \sum_{i=1}^m w_i \phi_i(x) \text{ with } \|\phi_i\|_1 \leq 1\}$$

- the convex hull of  $F$  is

$$\text{conv} F = \left\{ \sum_{i=1}^m \lambda_i f_i \mid f_i \in F; \lambda_i \geq 0; \sum_{i=1}^m \lambda_i = 1; m \in \mathbb{N} \right\}$$

- convex hull technique (Van Der Vaart et al., 1996, Theorem 2.6.9)

$$\log N_2(G_1; \|\cdot\|) \leq \log N_2(\bar{n}F; \|\cdot\|) \leq C \frac{1}{\bar{n}^{\frac{2d}{d+2}}}$$

- control the constant  $C$

$$C := \frac{D_k}{|\mathcal{Z}|} \left[ \frac{C_k}{|\mathcal{Z}|} (2^{d+1} + 1)^{\frac{1}{d}} \right]^{\frac{2d}{d+2}} \leq 10^7 d \quad \text{if } d > 5$$

# Proof sketch: convex hull technique and its constant!

- Consider the following function space

$$F = \{f : \mathbb{R}^d \rightarrow \mathbb{R} \mid f(x) = \sum_{i=1}^m w_i g_i(x) \mid w_i \in \mathbb{R}, g_i \in \mathcal{G}\}$$

- the convex hull of  $F$  is

$$\text{conv} F = \left\{ \sum_{i=1}^m \lambda_i f_i \mid f_i \in F, \lambda_i \geq 0, \sum_{i=1}^m \lambda_i = 1, m \in \mathbb{N} \right\}$$

- convex hull technique (Van Der Vaart et al., 1996, Theorem 2.6.9)

$$\log N_2(G; \|\cdot\|) \leq \log N_2(\text{conv} F; \|\cdot\|) \leq C \frac{1}{\epsilon}^{\frac{2d}{d+2}}$$

- control the constant  $C$

$$C := \frac{D_k}{\epsilon} \left[ \frac{C_k}{\epsilon} (2^{d+1} + 1)^{\frac{1}{d}} \right]^{\frac{2d}{d+2}} \leq 10^7 d \quad \text{if } d > 5$$

# Proof sketch: convex hull technique and its constant!

- Consider the following function space

$$F = \{f : \mathbb{R}^d \rightarrow \mathbb{R} \mid f(x) = \sum_{i=1}^m w_i \phi_i(x) \text{ with } \sum_{i=1}^m w_i^2 \leq 1\}$$

- the convex hull of  $F$  is

$$\text{conv} F = \left\{ \sum_{i=1}^m \lambda_i f_i \mid f_i \in F, \lambda_i \geq 0, \sum_{i=1}^m \lambda_i = 1, m \in \mathbb{N} \right\}$$

- convex hull technique (Van Der Vaart et al., 1996, Theorem 2.6.9)

$$\log N_2(G_1; \|\cdot\|) \leq \log N_2(\text{conv} F; \|\cdot\|) \leq C \int_0^1 \frac{1}{\sqrt{t}} dt$$

- control the constant  $C$

$$C := \frac{D_k}{\sqrt{2}} \left[ \frac{C_k}{\sqrt{2}} (2^{d+1} + 1)^{\frac{1}{d}} \right]^{\frac{2d}{d+2}} \leq 10^7 d \text{ if } d > 5$$

# Proof sketch: convex hull technique and its constant!

- Consider the following function space

$$F = \{f : \mathbb{R}^d \rightarrow \mathbb{R} \mid f(x) = \sum_{i=1}^m w_i g_i(x) \text{ with } w_i \geq 0, \sum_{i=1}^m w_i = 1\}$$

- the convex hull of  $F$  is

$$\text{conv} F = \left\{ \sum_{i=1}^m \lambda_i f_i \mid f_i \in F, \lambda_i \geq 0, \sum_{i=1}^m \lambda_i = 1, m \in \mathbb{N} \right\}$$

- convex hull technique (Van Der Vaart et al., 1996, Theorem 2.6.9)

$$\log N_2(G; \|\cdot\|) \leq \log N_2(\text{conv} F; \|\cdot\|) \leq C \frac{1}{d^{\frac{2d}{d+2}}}$$

- control the constant  $C$

$$C := \frac{D_k}{|\mathcal{Z}|} \left[ \frac{C_k}{|\mathcal{Z}|} (2^{d+1} + 1)^{\frac{1}{d}} \right]^{\frac{2d}{d+2}} \leq 10^7 d \quad \text{if } d > 5$$

= (d) = (1)