

**The role of over-parameterization in machine learning:
the good, the bad, the ugly**
- from a function space view

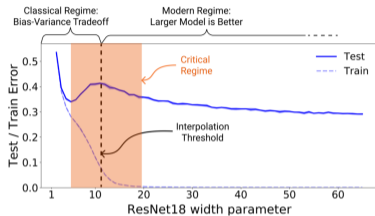
Fanghui Liu

Department of Computer Science, University of Warwick, UK
Centre for Discrete Mathematics and its Applications (DIMAP), Warwick

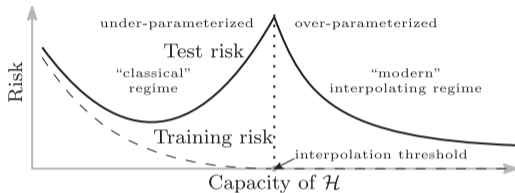
at AAI New Faculty Highlights 2024, Vancouver



Surprises in modern neural networks: double descent

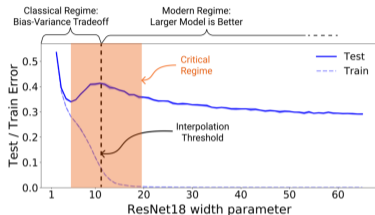


(a) Training and test error on ResNet18 [1]

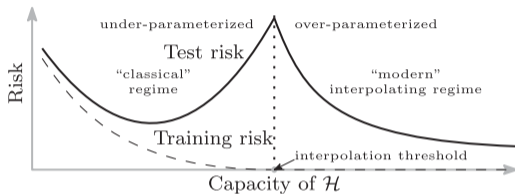


(b) Double descent [2] (Belkin, Hsu, Ma, Mandal, 2019).

Surprises in modern neural networks: double descent



(a) Training and test error on ResNet18 [1]

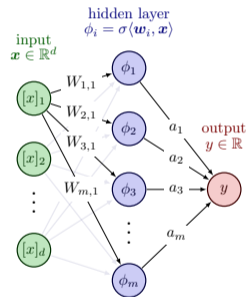
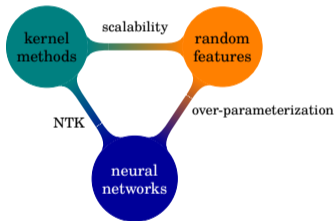
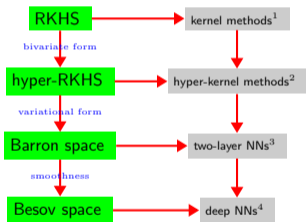


(b) Double descent [2] (Belkin, Hsu, Ma, Mandal, 2019).

Observations: beyond bias-variance trade-off

- ▶ 1) Monotonic decreasing in the overparameterized regime
- ▶ 2) Global minimum when #parameters is infinite
- ▶ 3) Peak at the interpolation thresholds

Today's talk: Function spaces vs. Models (initialization matters)



¹[LHGYL, JMLR20; LHCS, TPAMI21; LLS, AISTATS21]

²[LSHYS, JMLR21]

³[LSC, NeurIPS22; LHCS, TPAMI22; LHCS, AISTATS21]

⁴[LVC, NeurIPS22; ZLCC, NeurIPS22; WZLCC, NeurIPS22, ZLCLC, ICML23]

Questions on high dimensional kernel methods

- double descent based on random matrix theory: ([Mei and Montanari, 2022](#)), ([Hastie, Montanari, Rosset, Tibshirani, 2022](#)), ([Liao, Couillet, Mahoney, 2022](#))

Questions on high dimensional kernel methods

- double descent based on random matrix theory: ([Mei and Montanari, 2022](#)), ([Hastie, Montanari, Rosset, Tibshirani, 2022](#)), ([Liao, Couillet, Mahoney, 2022](#))

high dimensional kernel methods can only learn linear function! [3]

Questions on high dimensional kernel methods

- double descent based on random matrix theory: (Mei and Montanari, 2022), (Hastie, Montanari, Rosset, Tibshirani, 2022), (Liao, Couillet, Mahoney, 2022)

high dimensional kernel methods can only learn linear function! [3]

- asymptotic expansion under high dimensions [4] (El Karoui, 2010)
under the setting of $n, d \rightarrow \infty$, $n/d \rightarrow \psi_1$ as $d \rightarrow \infty$ with $\psi_1 \in (0, \infty)$, we have

$$\|\mathbf{K} - (a\mathbf{X}\mathbf{X}^\top + b\mathbf{I})\|_2 \xrightarrow{\mathbb{P}} 0 \quad \text{when } d \rightarrow \infty \quad \text{for some parameters } a, b$$

Questions on high dimensional kernel methods

- double descent based on random matrix theory: (Mei and Montanari, 2022), (Hastie, Montanari, Rosset, Tibshirani, 2022), (Liao, Couillet, Mahoney, 2022)

high dimensional kernel methods can only learn linear function! [3]

- asymptotic expansion under high dimensions [4] (El Karoui, 2010)
under the setting of $n, d \rightarrow \infty$, $n/d \rightarrow \psi_1$ as $d \rightarrow \infty$ with $\psi_1 \in (0, \infty)$, we have

$$\|\mathbf{K} - (a\mathbf{X}\mathbf{X}^\top + b\mathbf{I})\|_2 \xrightarrow{\mathbb{P}} 0 \quad \text{when } d \rightarrow \infty \quad \text{for some parameters } a, b$$

- $\|f\|_{\mathcal{H}} < \infty$?

Motivation

- ▶ high dimension vs. fixed dimension
- ▶ from asymptotic to non-asymptotic
- ▶ two-layer neural networks trained by SGD

Motivation

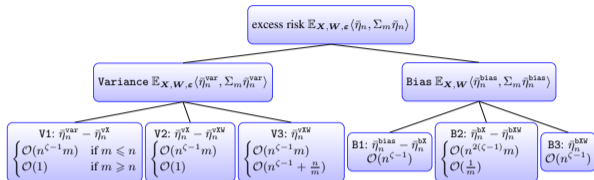
- ▶ high dimension vs. fixed dimension
- ▶ from asymptotic to non-asymptotic
- ▶ two-layer neural networks trained by SGD
- Analysis
 - ▶ dimension-free bound
 - ▶ multiple randomness sources
 - data sampling, label noise, Gaussian initialization, stochastic gradients

Motivation

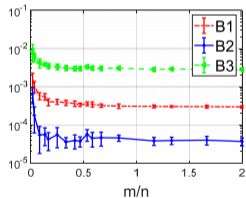
- ▶ high dimension vs. fixed dimension
- ▶ from asymptotic to non-asymptotic
- ▶ two-layer neural networks trained by SGD
- Analysis
 - ▶ dimension-free bound
 - ▶ multiple randomness sources
 - data sampling, label noise, Gaussian initialization, stochastic gradients

observations 1), 2), 3) can be still proved!

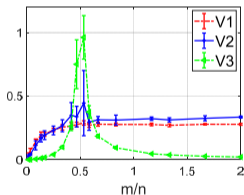
Our results: Double descent of RFMs trained by SGD¹



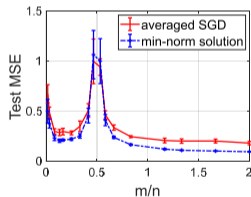
- ▶ (partially) decouple multiple randomness sources
- ▶ converge to $\mathcal{O}(1)$ order (noise variance)
- ▶ monotonic decreasing **bias** + unimodal **variance**
- ▶ constant step-size SGD does not hurt the convergence rate



(c) bias



(d) variance

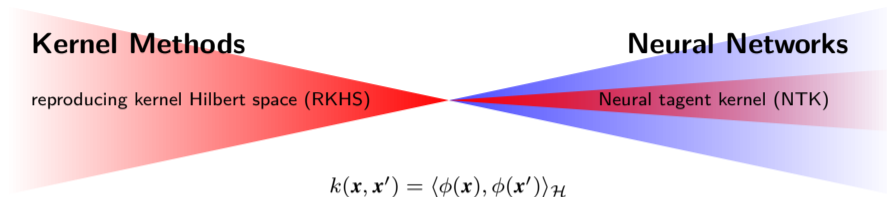


(e) excess risk

¹Fanghui Liu, Johan Suykens, Volkan Cevher. *On the Double Descent of Random Features Models Trained with SGD*. NeurIPS 2022.

Fanghui Liu, Xiaolin Huang, Yudong Chen, and Johan Suykens. *Random Features for Kernel Approximation: A Survey on Algorithms, Theory, and Beyond*. TPAMI2021.

From kernel methods (RKHS) to neural networks (?)



From kernel methods (RKHS) to neural networks (?)

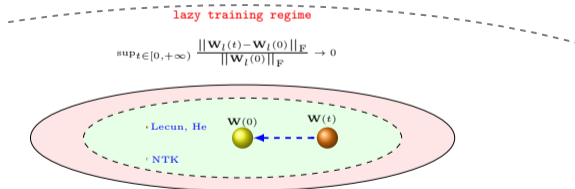
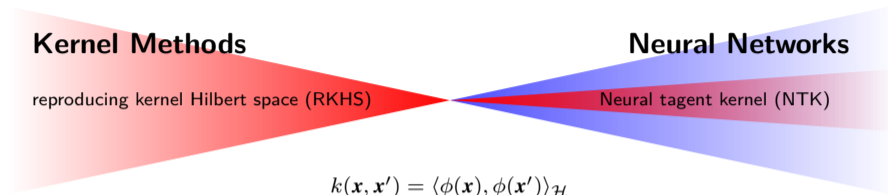
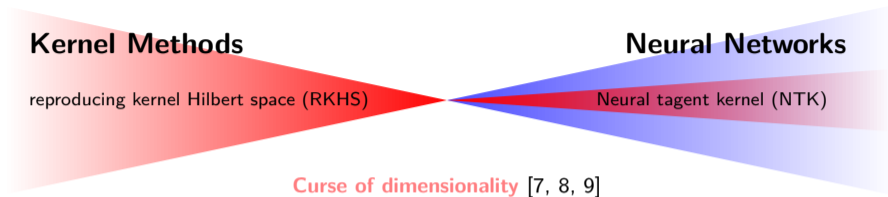


Figure: Lazy training regime: under the NTK initialization [5, 6].

From kernel methods (RKHS) to neural networks (?)



From kernel methods (RKHS) to neural networks (?)

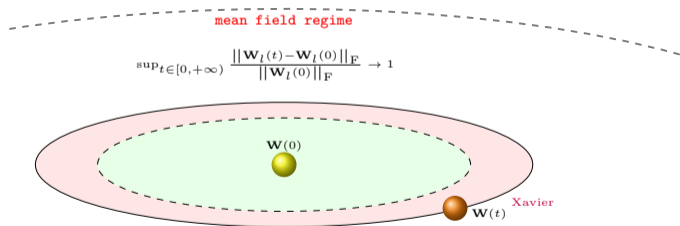
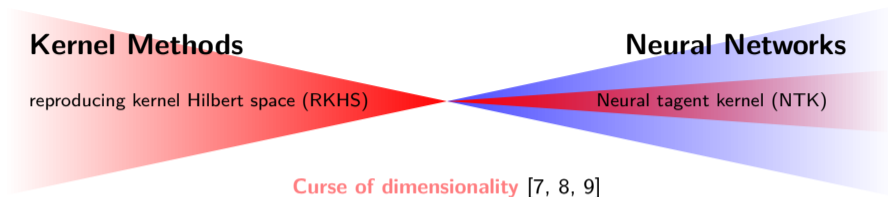


Figure: Mean field regime: under the Xavier initialization, abc-Parametrizations [10, 11].

From RKHS to Barron space

- o RKHS of RFMs:

$$\hat{k}_m(\mathbf{x}, \mathbf{x}') = \frac{1}{m} \sum_{i=1}^m \phi(\mathbf{x}, \mathbf{w}_i) \phi(\mathbf{x}', \mathbf{w}_i) \rightarrow k_\mu(\mathbf{x}, \mathbf{x}') = \int_{\mathcal{W}} \phi(\mathbf{x}, \mathbf{w}) \phi(\mathbf{x}', \mathbf{w}) d\mu(\mathbf{w})$$

From RKHS to Barron space

o RKHS of RFMs:

$$\hat{k}_m(\mathbf{x}, \mathbf{x}') = \frac{1}{m} \sum_{i=1}^m \phi(\mathbf{x}, \mathbf{w}_i) \phi(\mathbf{x}', \mathbf{w}_i) \rightarrow k_\mu(\mathbf{x}, \mathbf{x}') = \int_{\mathcal{W}} \phi(\mathbf{x}, \mathbf{w}) \phi(\mathbf{x}', \mathbf{w}) d\mu(\mathbf{w})$$

Definition (Barron space [12] (E, Ma, Wu, 2021))

$$\mathcal{B} = \cup_{\mu \in \mathcal{P}(\mathcal{W})} \mathcal{H}_{k_\mu}, \quad \|f\|_{\mathcal{B}} = \inf_{\mu \in \mathcal{P}(\mathcal{W})} \|f\|_{\mathcal{H}_{k_\mu}}$$

From RKHS to Barron space

o RKHS of RFMs:

$$\hat{k}_m(\mathbf{x}, \mathbf{x}') = \frac{1}{m} \sum_{i=1}^m \phi(\mathbf{x}, \mathbf{w}_i) \phi(\mathbf{x}', \mathbf{w}_i) \rightarrow k_\mu(\mathbf{x}, \mathbf{x}') = \int_{\mathcal{W}} \phi(\mathbf{x}, \mathbf{w}) \phi(\mathbf{x}', \mathbf{w}) d\mu(\mathbf{w})$$

Definition (Barron space [12] (E, Ma, Wu, 2021))

$$\mathcal{B} = \cup_{\mu \in \mathcal{P}(\mathcal{W})} \mathcal{H}_{k_\mu}, \quad \|f\|_{\mathcal{B}} = \inf_{\mu \in \mathcal{P}(\mathcal{W})} \|f\|_{\mathcal{H}_{k_\mu}}$$

Remark: o Two-layer neural networks: data-adaptive kernel

From RKHS to Barron space

◦ RKHS of RFMs:

$$\hat{k}_m(\mathbf{x}, \mathbf{x}') = \frac{1}{m} \sum_{i=1}^m \phi(\mathbf{x}, \mathbf{w}_i) \phi(\mathbf{x}', \mathbf{w}_i) \rightarrow k_\mu(\mathbf{x}, \mathbf{x}') = \int_{\mathcal{W}} \phi(\mathbf{x}, \mathbf{w}) \phi(\mathbf{x}', \mathbf{w}) d\mu(\mathbf{w})$$

Definition (Barron space [12] (E, Ma, Wu, 2021))

$$\mathcal{B} = \cup_{\mu \in \mathcal{P}(\mathcal{W})} \mathcal{H}_{k_\mu}, \quad \|f\|_{\mathcal{B}} = \inf_{\mu \in \mathcal{P}(\mathcal{W})} \|f\|_{\mathcal{H}_{k_\mu}}$$

Remark: ◦ Two-layer neural networks: data-adaptive kernel

◦ equivalent to path norm $\|\Theta\|_{\mathcal{P}} := \frac{1}{m} \sum_{k=1}^m |a_k| \|\mathbf{w}_k\|_1$

From RKHS to Barron space

- o RKHS of RFMs:

$$\hat{k}_m(\mathbf{x}, \mathbf{x}') = \frac{1}{m} \sum_{i=1}^m \phi(\mathbf{x}, \mathbf{w}_i) \phi(\mathbf{x}', \mathbf{w}_i) \rightarrow k_\mu(\mathbf{x}, \mathbf{x}') = \int_{\mathcal{W}} \phi(\mathbf{x}, \mathbf{w}) \phi(\mathbf{x}', \mathbf{w}) d\mu(\mathbf{w})$$

Definition (Barron space [12] (E, Ma, Wu, 2021))

$$\mathcal{B} = \cup_{\mu \in \mathcal{P}(\mathcal{W})} \mathcal{H}_{k_\mu}, \quad \|f\|_{\mathcal{B}} = \inf_{\mu \in \mathcal{P}(\mathcal{W})} \|f\|_{\mathcal{H}_{k_\mu}}$$

Remark: o Two-layer neural networks: data-adaptive kernel

o equivalent to path norm $\|\Theta\|_{\mathcal{P}} := \frac{1}{m} \sum_{k=1}^m |a_k| \|\mathbf{w}_k\|_1$

o parameter space vs. measure space

e.g., [7] (Bach, 2017), [13] (Bartolucci, Vito, Rosasco, Vigogna, 2022).

Our results: Refined analyses in Barron spaces²

For the class of two-layer neural networks \mathcal{F}_m

$$\theta^* = \arg \min_{f_{\theta} \in \mathcal{F}_m} \frac{1}{n} \sum_{i=1}^n (y_i - f_{\theta}(\mathbf{x}_i))^2 + \lambda \|\theta\|_{\mathcal{P}}.$$

²Fanghui Liu, Leello Dadi, Volkan Cevher. Learning with two-layer, norm-constrained, over-parameterized neural networks. JMLR ([under the second-round review](#))

Our results: Refined analyses in Barron spaces²

For the class of two-layer neural networks \mathcal{F}_m

$$\theta^* = \arg \min_{f_{\theta} \in \mathcal{F}_m} \frac{1}{n} \sum_{i=1}^n (y_i - f_{\theta}(\mathbf{x}_i))^2 + \lambda \|\theta\|_{\mathcal{P}}.$$

Theorem (Informal)

Under proper assumptions, for two-layer *over-parameterized* neural networks, learning in Barron spaces leads to

$$\|f_{\theta^*} - f_{\rho}\|_{L_{\rho_X}^2}^2 \lesssim \lambda + \frac{1}{m} + d^2 n^{-\frac{d+2}{2d+2}} \quad w.h.p.$$

²Fanghui Liu, Leello Dadi, Volkan Cevher. Learning with two-layer, norm-constrained, over-parameterized neural networks. JMLR ([under the second-round review](#))

Our results: Refined analyses in Barron spaces²

For the class of two-layer neural networks \mathcal{F}_m

$$\theta^* = \arg \min_{f_\theta \in \mathcal{F}_m} \frac{1}{n} \sum_{i=1}^n (y_i - f_\theta(\mathbf{x}_i))^2 + \lambda \|\theta\|_{\mathcal{P}}.$$

Theorem (Informal)

Under proper assumptions, for two-layer *over-parameterized* neural networks, learning in Barron spaces leads to

$$\|f_{\theta^*} - f_\rho\|_{L^2_{\rho_X}}^2 \lesssim \lambda + \frac{1}{m} + d^2 n^{-\frac{d+2}{2d+2}} \quad w.h.p.$$

Remark:

- ▶ [14] (Siegel, Xu, 2022) on metric entropy

$$\epsilon^{-\frac{2d}{d+3}} d \lesssim \log \mathcal{N}_2(\mathcal{G}_1, \epsilon) \lesssim d \epsilon^{-\frac{2d}{d+3}}.$$

²Fanghui Liu, Leello Dadi, Volkan Cevher. Learning with two-layer, norm-constrained, over-parameterized neural networks. JMLR ([under the second-round review](#))

Our results: Refined analyses in Barron spaces²

For the class of two-layer neural networks \mathcal{F}_m

$$\theta^* = \arg \min_{f_\theta \in \mathcal{F}_m} \frac{1}{n} \sum_{i=1}^n (y_i - f_\theta(\mathbf{x}_i))^2 + \lambda \|\theta\|_{\mathcal{P}}.$$

Theorem (Informal)

Under proper assumptions, for two-layer *over-parameterized* neural networks, learning in Barron spaces leads to

$$\|f_{\theta^*} - f_\rho\|_{L^2_{\rho_X}}^2 \lesssim \lambda + \frac{1}{m} + d^2 n^{-\frac{d+2}{2d+2}} \quad w.h.p.$$

Remark:

- ▶ [14] (Siegel, Xu, 2022) on metric entropy

$$\epsilon^{-\frac{2d}{d+3}} d \lesssim \log \mathcal{N}_2(\mathcal{G}_1, \epsilon) \lesssim d \epsilon^{-\frac{2d}{d+3}} \leq 6144d^5 \epsilon^{-\frac{2d}{d+2}} \quad \text{[Ours]}$$

²Fanghui Liu, Leello Dadi, Volkan Cevher. Learning with two-layer, norm-constrained, over-parameterized neural networks. JMLR (under the second-round review)

Optimization in Barron spaces is difficult: curse of dimensionality!



Optimization in Barron spaces is difficult: curse of dimensionality!



What is the suitable function space of NNs, both **statistically** and **computationally** efficient?

Applications: Over-parameterization helps/hurts robustness?³

Helps! [15]



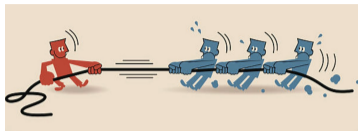
Hurts! [16, 17, 18]

³Zhenyu Zhu, **Fanghui Liu**, Grigorios Chrysos, Volkan Cevher. *Robustness in deep learning: The good (width), the bad (depth), and the ugly (initialization)*. NeurIPS 2022.

Jiayuan Ye, Zhenyu Zhu, **Fanghui Liu**, Reza Shokri, Volkan Cevher. Initialization matters: Privacy-utility analysis of overparameterized neural networks. NeurIPS 2023.

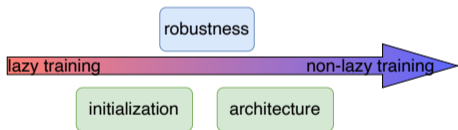
Applications: Over-parameterization helps/hurts robustness?³

Helps! [15]



Hurts! [16, 17, 18]

- ▶ initialization (e.g., lazy training, non-lazy training)
- ▶ architecture (e.g., width, depth)

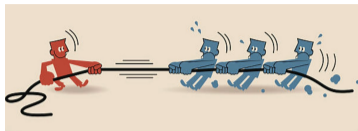


³Zhenyu Zhu, **Fanghui Liu**, Grigorios Chrysos, Volkan Cevher. *Robustness in deep learning: The good (width), the bad (depth), and the ugly (initialization)*. NeurIPS 2022.

Jiayuan Ye, Zhenyu Zhu, **Fanghui Liu**, Reza Shokri, Volkan Cevher. Initialization matters: Privacy-utility analysis of overparameterized neural networks. NeurIPS 2023.

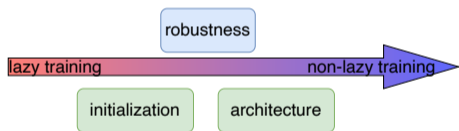
Applications: Over-parameterization helps/hurts robustness?³

Helps! [15]



Hurts! [16, 17, 18]

- ▶ initialization (e.g., lazy training, non-lazy training)
- ▶ architecture (e.g., width, depth)



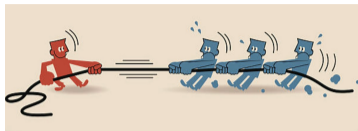
Takeaway messages: **the good (width), the bad (depth), the ugly (initialization)**

³Zhenyu Zhu, **Fanghui Liu**, Grigorios Chrysos, Volkan Cevher. *Robustness in deep learning: The good (width), the bad (depth), and the ugly (initialization)*. NeurIPS 2022.

Jiayuan Ye, Zhenyu Zhu, **Fanghui Liu**, Reza Shokri, Volkan Cevher. Initialization matters: Privacy-utility analysis of overparameterized neural networks. NeurIPS 2023.

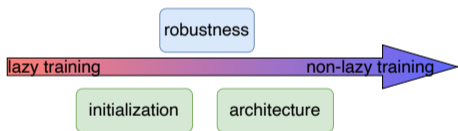
Applications: Over-parameterization helps/hurts robustness?³

Helps! [15]



Hurts! [16, 17, 18]

- ▶ initialization (e.g., lazy training, non-lazy training)
- ▶ architecture (e.g., width, depth)



Takeaway messages: **the good (width), the bad (depth), the ugly (initialization)**

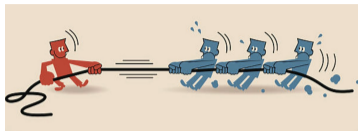
- ▶ width **helps** robustness in the over-parameterized regime

³Zhenyu Zhu, **Fanghui Liu**, Grigorios Chrysos, Volkan Cevher. *Robustness in deep learning: The good (width), the bad (depth), and the ugly (initialization)*. NeurIPS 2022.

Jiayuan Ye, Zhenyu Zhu, **Fanghui Liu**, Reza Shokri, Volkan Cevher. Initialization matters: Privacy-utility analysis of overparameterized neural networks. NeurIPS 2023.

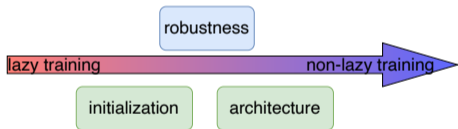
Applications: Over-parameterization helps/hurts robustness?³

Helps! [15]



Hurts! [16, 17, 18]

- ▶ initialization (e.g., lazy training, non-lazy training)
- ▶ architecture (e.g., width, depth)



Takeaway messages: **the good (width), the bad (depth), the ugly (initialization)**

- ▶ width **helps** robustness in the over-parameterized regime
- ▶ depth **helps** robustness in LeCun initialization but **hurts** robustness in He/NTK initialization

³Zhenyu Zhu, **Fanghui Liu**, Grigorios Chrysos, Volkan Cevher. *Robustness in deep learning: The good (width), the bad (depth), and the ugly (initialization)*. NeurIPS 2022.

Jiayuan Ye, Zhenyu Zhu, **Fanghui Liu**, Reza Shokri, Volkan Cevher. Initialization matters: Privacy-utility analysis of overparameterized neural networks. NeurIPS 2023.

Conclusion: the good, the bad, the ugly



	good	bad	ugly
kernel methods	analysis	performance	curse of dimensionality
neural networks	performance	analysis	over-parameterization
generalization	benign overfitting	catastrophic overfitting	model complexity
robustness	width	depth	initialization
privacy	depth	width	initialization

- ▶ IEEE ICASSP 2023 Tutorial - “Neural networks: the good, the bad, and the ugly”
- ▶ CVPR 2023 Tutorial - “Deep learning theory for computer vision”

Thanks for your attention!

Q & A

my homepage www.lfhsgre.org for more information!

References I

- [1] 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.
(Cited on pages 3 and 4.)
- [2] 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.
(Cited on pages 3 and 4.)
- [3] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Linearized two-layers neural networks in high dimension. *Annals of Statistics*, 49(2):1029–1054, 2021.
(Cited on pages 6, 7, 8, and 9.)
- [4] Nouredine El Karoui. The spectrum of kernel random matrices. *Annals of Statistics*, 38(1):1–50, 2010.
(Cited on pages 6, 7, 8, and 9.)

References II

- [5] Arthur Jacot, Franck Gabriel, and Clément Hongler.
Neural tangent kernel: Convergence and generalization in neural networks.
In Advances in Neural Information Processing Systems, pages 8571–8580, 2018.
(Cited on pages 14 and 15.)
- [6] Lenaïc Chizat, Edouard Oyallon, and Francis Bach.
On lazy training in differentiable programming.
In Advances in Neural Information Processing Systems, pages 2933–2943, 2019.
(Cited on pages 14 and 15.)
- [7] Francis Bach.
Breaking the curse of dimensionality with convex neural networks.
Journal of Machine Learning Research, 18(1):629–681, 2017.
(Cited on pages 16, 17, 18, 19, 20, 21, and 22.)
- [8] Gilad Yehudai and Ohad Shamir.
On the power and limitations of random features for understanding neural networks.
In Advances in Neural Information Processing Systems, pages 6594–6604, 2019.
(Cited on pages 16 and 17.)

References III

- [9] Michael Celentano, Theodor Misiakiewicz, and Andrea Montanari.
Minimum complexity interpolation in random features models.
arXiv preprint arXiv:2103.15996, 2021.
(Cited on pages 16 and 17.)
- [10] Greg Yang and Edward J Hu.
Feature learning in infinite-width neural networks.
arXiv preprint arXiv:2011.14522, 2020.
(Cited on pages 16 and 17.)
- [11] Lénaïc Chizat and Francis Bach.
On the global convergence of gradient descent for over-parameterized models using optimal transport.
Advances in Neural Information Processing Systems, 31:3036–3046, 2018.
(Cited on pages 16 and 17.)
- [12] 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.
(Cited on pages 18, 19, 20, 21, and 22.)

References IV

- [13] Francesca Bartolucci, Ernesto De Vito, Lorenzo Rosasco, and Stefano Vigogna. Understanding neural networks with reproducing kernel Banach spaces. *Applied and Computational Harmonic Analysis*, 2022.
(Cited on pages 18, 19, 20, 21, and 22.)
- [14] 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.
(Cited on pages 23, 24, 25, and 26.)
- [15] Sébastien Bubeck and Mark Sellke. A universal law of robustness via isoperimetry. In *Advances in Neural Information Processing Systems*, pages 28811–28822, 2021.
(Cited on pages 29, 30, 31, 32, and 33.)
- [16] Hamed Hassani and Adel Javanmard. The curse of overparametrization in adversarial training: Precise analysis of robust generalization for random features regression. *arXiv preprint arXiv:2201.05149*, 2022.
(Cited on pages 29, 30, 31, 32, and 33.)

References V

- [17] Boxi Wu, Jinghui Chen, Deng Cai, Xiaofei He, and Quanquan Gu.
Do wider neural networks really help adversarial robustness?
In *Advances in Neural Information Processing Systems*, pages 7054–7067, 2021.
(Cited on pages 29, 30, 31, 32, and 33.)
- [18] Hanxun Huang, Yisen Wang, Sarah Erfani, Quanquan Gu, James Bailey, and Xingjun Ma.
Exploring architectural ingredients of adversarially robust deep neural networks.
In *Advances in Neural Information Processing Systems*, pages 5545–5559, 2021.
(Cited on pages 29, 30, 31, 32, and 33.)