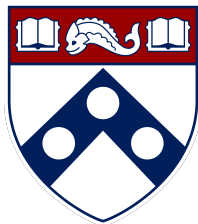


Non-Linear Feature Learning with One Gradient Step in Two-Layer Neural Networks

Behrad Moniri

University of Pennsylvania
bemoniri@seas.upenn.edu

March 23rd, 2024





Collaborators



Donghwan Lee



Hamed Hassani



Edgar Dobriban



Introduction



Deep Learning is *very* successful.



Introduction

- It is a huge *engineering* success.



Introduction

- It is a huge *engineering* success.
- Why is it very successful?



Introduction

- It is a huge *engineering* success.
- Why is it very successful?
- The problem is notoriously hard.



- Is there a *rich enough*, solvable model?



- Is there a *rich enough*, solvable model?
- There are **too** many moving parts: architecture, optimization, data, etc.



We consider the simplest possible architecture.



We consider the simplest possible architecture.

A two-layer fully connected neural network.



Two-Layer Neural Networks

Input
 $\mathbf{x} \in \mathbb{R}^d$

$x[1]$

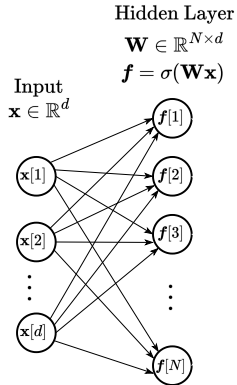
$x[2]$

\vdots

$x[d]$

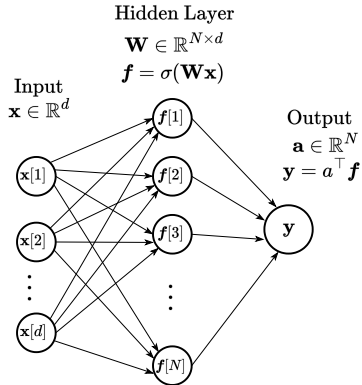


Two-Layer Neural Networks



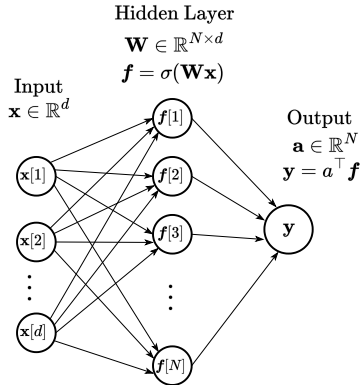


Two-Layer Neural Networks

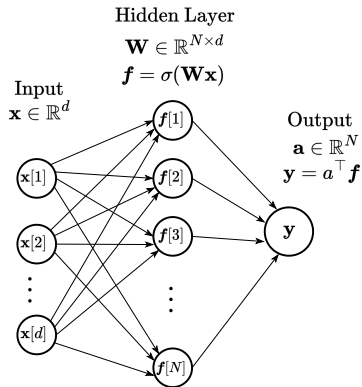




Two-Layer Neural Networks



Asymptotic Regime: $n, d, N \rightarrow \infty$ with $d/n \rightarrow \phi$ and $d/N \rightarrow \psi$.



- **Simplest Model:**
Random Features Model.
(Rahimi and Recht, 2007)



Random Features Models

- Random features model is **very popular**:



- Random features model is **very popular**:
 - Used to study various aspects of deep learning such as double descent, robustness to adversarial attacks, privacy, fairness, OOD performance, calibration, etc.

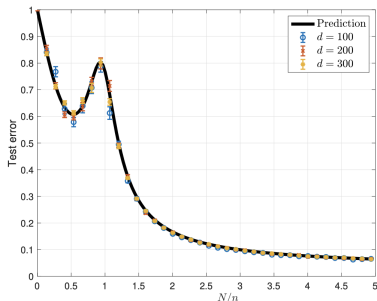
See e.g., Mei and Montanari (2022); Lin and Dobriban (2021); Lee et al. (2023); Hassani and Javanmard (2022); Bombari and Mondelli (2023); Bombari et al. (2023); Clarté et al. (2023), etc.



Random Features Models

- Random features model is **very popular**:
 - Used to study various aspects of deep learning such as double descent, robustness to adversarial attacks, privacy, fairness, OOD performance, calibration, etc.

See e.g., Mei and Montanari (2022); Lin and Dobriban (2021); Lee et al. (2023); Hassani and Javanmard (2022); Bombari and Mondelli (2023); Bombari et al. (2023); Clarté et al. (2023), etc.



Double descent in random feature models (Mei and Montanari, 2022).



Random Features Models

- Random features models can only learn linear functions.

Mei and Montanari (2022); Hu and Lu (2023)



Random Features Models

- Random features models can only learn linear functions.

Mei and Montanari (2022); Hu and Lu (2023)

- **Gaussian Equivalence:**



- Random features models can only learn linear functions.

Mei and Montanari (2022); Hu and Lu (2023)

- **Gaussian Equivalence:**

$$\begin{aligned}\sigma(\mathbf{W}\mathbf{x}) &= c_1\mathbf{W}\mathbf{x} + c_2H_2(\mathbf{W}\mathbf{x}) + \dots \\ &\approx c_1\mathbf{W}\mathbf{x} + \mathbf{z}\end{aligned}$$



- Random features models can only learn linear functions.

Mei and Montanari (2022); Hu and Lu (2023)

- **Gaussian Equivalence:**

$$\begin{aligned}\sigma(\mathbf{W}\mathbf{x}) &= c_1\mathbf{W}\mathbf{x} + c_2H_2(\mathbf{W}\mathbf{x}) + \dots \\ &\approx c_1\mathbf{W}\mathbf{x} + \mathbf{z}\end{aligned}$$

- This makes analysis *easy* but the model very *limited*.



- Random features models can only learn linear functions.

Mei and Montanari (2022); Hu and Lu (2023)

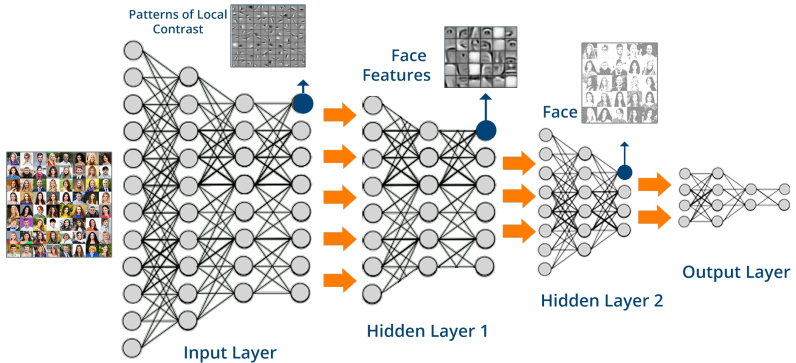
- **Gaussian Equivalence:**

$$\begin{aligned}\sigma(\mathbf{W}\mathbf{x}) &= c_1\mathbf{W}\mathbf{x} + c_2H_2(\mathbf{W}\mathbf{x}) + \dots \\ &\approx c_1\mathbf{W}\mathbf{x} + \mathbf{z}\end{aligned}$$

- This makes analysis *easy* but the model very *limited*.
- What is missing?



Feature Learning





- Feature learning is absent in random feature models.



Random Features Models

- Feature learning is absent in random feature models.
- How to go beyond random feature models?



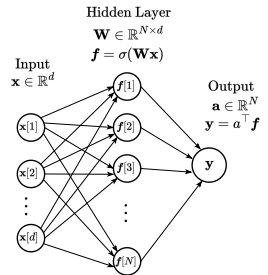
Random Features Models

- Feature learning is absent in random feature models.
- How to go beyond random feature models?
- Let's do one step of gradient descent on first layer weights.

One Gradient Step



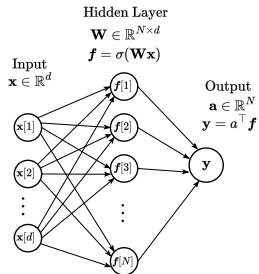
One Gradient Step Update





One Gradient Step Update

We train the network as follows:



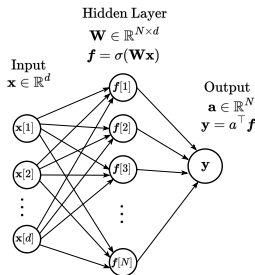


One Gradient Step Update

We train the network as follows:

1. We first initialize

$$\mathbf{a} \sim \mathcal{N}\left(\mathbf{0}_N, \frac{1}{N}\mathbf{I}_N\right), \quad \text{and} \quad [\mathbf{W}_0]_{ij} \sim \mathcal{N}\left(0, \frac{1}{d}\right)$$





One Gradient Step Update

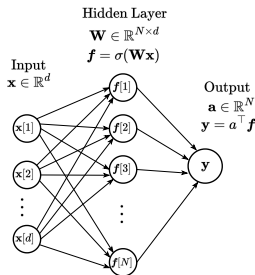
We train the network as follows:

1. We first initialize

$$\mathbf{a} \sim \mathcal{N}\left(\mathbf{0}_N, \frac{1}{N}\mathbf{I}_N\right), \quad \text{and} \quad [\mathbf{W}_0]_{ij} \sim \mathcal{N}\left(0, \frac{1}{d}\right)$$

2. We take one gradient step on the empirical MSE loss

$$\mathbf{W}_1 = \mathbf{W}_0 - \eta \frac{\partial}{\partial \mathbf{W}} \left(\|\mathbf{y} - \sigma(\mathbf{X}\mathbf{W}^\top)\mathbf{a}\|_2^2 \right) \Big|_{\mathbf{W}_0, \mathbf{a}}$$





One Gradient Step Update

We train the network as follows:

1. We first initialize

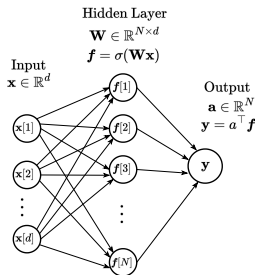
$$\mathbf{a} \sim \mathcal{N}\left(\mathbf{0}_N, \frac{1}{N}\mathbf{I}_N\right), \quad \text{and} \quad [\mathbf{W}_0]_{ij} \sim \mathcal{N}\left(0, \frac{1}{d}\right)$$

2. We take one gradient step on the empirical MSE loss

$$\mathbf{W}_1 = \mathbf{W}_0 - \eta \frac{\partial}{\partial \mathbf{W}} \left(\|\mathbf{y} - \sigma(\mathbf{X}\mathbf{W}^\top)\mathbf{a}\|_2^2 \right) \Big|_{\mathbf{W}_0, \mathbf{a}}$$

3. Fit \mathbf{a} via ridge regression:

$$\hat{\mathbf{a}} = \arg \min_{\mathbf{a} \in \mathbb{R}^N} \frac{1}{n} \|\mathbf{y} - \mathbf{F}\mathbf{a}\|_2^2 + \lambda \|\mathbf{a}\|_2^2, \quad \mathbf{F} = \sigma(\mathbf{X}\mathbf{W}_1^\top) \in \mathbb{R}^{n \times N}.$$





- **Data generation:**

$$x_i \stackrel{i.i.d.}{\sim} \mathbf{N}(0, \mathbf{I}_d), \quad y_i = f_{\star}(\mathbf{x}_i) + \varepsilon_i.$$



- **Data generation:**

$$\mathbf{x}_i \stackrel{i.i.d.}{\sim} \mathbf{N}(0, \mathbf{I}_d), \quad y_i = f_\star(\mathbf{x}_i) + \varepsilon_i.$$

- With one-step, only a single-index approximation can be learned. Thus, we let

$$f_\star(\mathbf{x}_i) = \sigma_\star(\boldsymbol{\beta}_\star^\top \mathbf{x}_i)$$

(see e.g., Dandi et al. (2023), etc.)



Partial answers in the prior work:



Partial answers in the prior work:

- Ba et al. (2022) show that if $\eta = O(1)$, still no nonlinear component of the teacher function can be learned.



Partial answers in the prior work:

- Ba et al. (2022) show that if $\eta = O(1)$, still no nonlinear component of the teacher function can be learned.
- Performance is still worse than linear regression.



Prior Work



- However, with $\eta = O(\sqrt{n})$ nonlinear functions can be learned.



- However, with $\eta = O(\sqrt{n})$ nonlinear functions can be learned.
- How is this possible? What happens?

Spectral Analysis



Spectrum of the Weights

- After the update, the weights are

$$\begin{aligned}\mathbf{W}_1 &= \mathbf{W}_0 - \eta \frac{\partial}{\partial \mathbf{W}} \left(\|\mathbf{y} - \sigma(\mathbf{X}\mathbf{W}^\top)\mathbf{a}\|_2^2 \right) \Big|_{\mathbf{W}_0, \mathbf{a}} \\ &= \mathbf{W}_0 + \frac{\eta}{n} \left[(\mathbf{a}\mathbf{y}^\top - \mathbf{a}\mathbf{a}^\top \sigma(\mathbf{W}_0\mathbf{X}^\top)) \circ \sigma'(\mathbf{W}_0\mathbf{X}^\top) \right] \mathbf{X},\end{aligned}$$



Spectrum of the Weights

- After the update, the weights are

$$\begin{aligned}\mathbf{W}_1 &= \mathbf{W}_0 - \eta \frac{\partial}{\partial \mathbf{W}} \left(\|\mathbf{y} - \sigma(\mathbf{X}\mathbf{W}^\top)\mathbf{a}\|_2^2 \right) \Big|_{\mathbf{W}_0, \mathbf{a}} \\ &= \mathbf{W}_0 + \frac{\eta}{n} \left[(\mathbf{a}\mathbf{y}^\top - \mathbf{a}\mathbf{a}^\top \sigma(\mathbf{W}_0\mathbf{X}^\top)) \circ \sigma'(\mathbf{W}_0\mathbf{X}^\top) \right] \mathbf{X},\end{aligned}$$

- Orthogonal decomposition:

$$\sigma'(\mathbf{W}_0\mathbf{X}^\top) = c_1 + \sigma'_\perp(\mathbf{W}_0\mathbf{X}^\top), \quad \text{with} \quad \mathbb{E}\sigma'_\perp(\mathbf{W}_0\mathbf{X}^\top) = \mathbf{0}$$



Spectrum of the Weights

- After the update, the weights are

$$\begin{aligned}\mathbf{W}_1 &= \mathbf{W}_0 - \eta \frac{\partial}{\partial \mathbf{W}} \left(\|\mathbf{y} - \sigma(\mathbf{X}\mathbf{W}^\top)\mathbf{a}\|_2^2 \right) \Big|_{\mathbf{W}_0, \mathbf{a}} \\ &= \mathbf{W}_0 + \frac{\eta}{n} \left[(\mathbf{a}\mathbf{y}^\top - \mathbf{a}\mathbf{a}^\top \sigma(\mathbf{W}_0\mathbf{X}^\top)) \circ \sigma'(\mathbf{W}_0\mathbf{X}^\top) \right] \mathbf{X},\end{aligned}$$

- Orthogonal decomposition:

$$\sigma'(\mathbf{W}_0\mathbf{X}^\top) = c_1 + \sigma'_\perp(\mathbf{W}_0\mathbf{X}^\top), \quad \text{with} \quad \mathbb{E}\sigma'_\perp(\mathbf{W}_0\mathbf{X}^\top) = \mathbf{0}$$

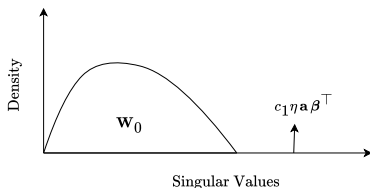
- Rank-1 Approximation:

$$\mathbf{W}_1 = \mathbf{W}_0 + \eta c_1 \mathbf{a} \left(\frac{\mathbf{X}^\top \mathbf{y}}{n} \right)^\top + \text{small.}$$



Spectrum of the Weights

$$\mathbf{W}_1 = \mathbf{W}_0 + \eta c_1 \mathbf{a} \left(\frac{\mathbf{X}^\top \mathbf{y}}{n} \right)^\top + \text{small.}$$



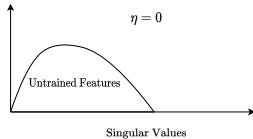
The vector $\beta := \frac{\mathbf{X}^\top \mathbf{y}}{n}$ is aligned to β_* .



Spectrum of the Feature Matrix

Updated Weight Matrix: $\mathbf{W} = \mathbf{W}_0 + \eta c_1 \mathbf{a} \boldsymbol{\beta}^\top$

Feature Matrix: $\mathbf{F} = \sigma(\mathbf{XW}^\top) = \sigma(\mathbf{XW}_0^\top + c_1 \eta \mathbf{X} \boldsymbol{\beta} \mathbf{a}^\top) \in \mathbb{R}^{n \times N}$

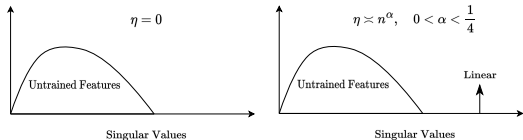




Spectrum of the Feature Matrix

Updated Weight Matrix: $\mathbf{W} = \mathbf{W}_0 + \eta c_1 \mathbf{a} \boldsymbol{\beta}^\top$

Feature Matrix: $\mathbf{F} = \sigma(\mathbf{XW}^\top) = \sigma(\mathbf{XW}_0^\top + c_1 \eta \mathbf{X} \boldsymbol{\beta} \mathbf{a}^\top) \in \mathbb{R}^{n \times N}$

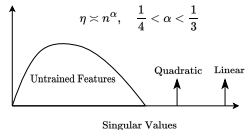
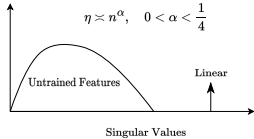
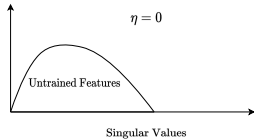




Spectrum of the Feature Matrix

Updated Weight Matrix: $\mathbf{W} = \mathbf{W}_0 + \eta c_1 \mathbf{a} \boldsymbol{\beta}^\top$

Feature Matrix: $\mathbf{F} = \sigma(\mathbf{XW}^\top) = \sigma(\mathbf{XW}_0^\top + c_1 \eta \mathbf{X} \boldsymbol{\beta} \mathbf{a}^\top) \in \mathbb{R}^{n \times N}$





Theorem

Let $\eta \asymp n^\alpha$ with $\frac{\ell-1}{2\ell} < \alpha < \frac{\ell}{2\ell+2}$ for some $\ell \in \mathbb{N}$. We have

$$\mathbf{F} = \mathbf{F}_\ell + \mathbf{\Delta}, \text{ with } \mathbf{F}_\ell := \mathbf{F}_0 + \sum_{k=1}^{\ell} c_1^k c_k \eta^k (\mathbf{X}\boldsymbol{\beta})^{\circ k} (\mathbf{a}^{\circ k})^\top,$$

where $\|\mathbf{\Delta}\|_{\text{op}} = o(\sqrt{n})$ with probability $1 - o(1)$.



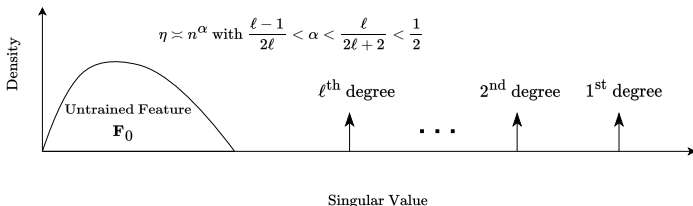
Spectrum of Updated Feature Matrix

Theorem

Let $\eta \asymp n^\alpha$ with $\frac{\ell-1}{2\ell} < \alpha < \frac{\ell}{2\ell+2}$ for some $\ell \in \mathbb{N}$. We have

$$\mathbf{F} = \mathbf{F}_\ell + \mathbf{\Delta}, \text{ with } \mathbf{F}_\ell := \mathbf{F}_0 + \sum_{k=1}^{\ell} c_1^k c_k \eta^k (\mathbf{X}\boldsymbol{\beta})^{\circ k} (\mathbf{a}^{\circ k})^\top,$$

where $\|\mathbf{\Delta}\|_{\text{op}} = o(\sqrt{n})$ with probability $1 - o(1)$.



Analysis of the Training/Test Error



- What error does the trained neural network achieve?



Universality

- What error does the trained neural network achieve?
- To do this, we show that the limiting behavior of test/train error is unchanged if



- What error does the trained neural network achieve?
- To do this, we show that the limiting behavior of test/train error is unchanged if
 - 1 we replace $\mathbf{F} = \sigma(\mathbf{X}\mathbf{W}_1^\top)$ with

$$\mathbf{F} = \mathbf{F}_0 + \sum_{k=1}^{\ell} c_1^k c_k \eta^k(\mathbf{X}\boldsymbol{\beta})^{ok} (\mathbf{a}^{ok})^\top.$$



- What error does the trained neural network achieve?
- To do this, we show that the limiting behavior of test/train error is unchanged if

1 we replace $\mathbf{F} = \sigma(\mathbf{X}\mathbf{W}_1^\top)$ with

$$\mathbf{F} = \mathbf{F}_0 + \sum_{k=1}^{\ell} c_1^k c_k \eta^k(\mathbf{X}\boldsymbol{\beta})^{ok} (\mathbf{a}^{ok})^\top.$$

2 we replace \mathbf{F}_0 with $\mathbf{F}_0 = c_1 \mathbf{X}\mathbf{W}_0^\top + c_{>1} \mathbf{Z}$.



Analysis of the Training/Test Errors

With these universality results, we can analyze the limiting behaviour of the train/test errors.



With these universality results, we can analyze the limiting behaviour of the train/test errors.

Theorem

Let $\ell \in \mathbb{N}$ and $\eta \asymp n^\alpha$ with $\frac{\ell-1}{2\ell} < \alpha < \frac{\ell}{2\ell+2}$, then for the learned feature map \mathbf{F} and the untrained feature map \mathbf{F}_0 , we have

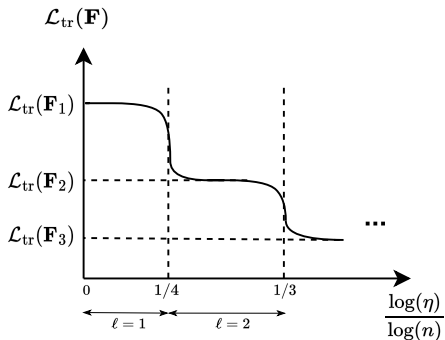
$$\mathcal{L}_{\text{tr}}(\mathbf{F}_0) - \mathcal{L}_{\text{tr}}(\mathbf{F}) \rightarrow_P \Delta_{\text{tr}} > 0$$

$$\mathcal{L}_{\text{te}}(\mathbf{F}_0) - \mathcal{L}_{\text{te}}(\mathbf{F}) \rightarrow_P \Delta_{\text{te}} > 0$$

The expression for $\Delta_{\text{te/tr}}$ can be found in the paper.



Analysis of the Training/Test Errors



Simulations



We consider

$$\textbf{Setting 1} : y = H_1(\boldsymbol{\beta}_*^\top \mathbf{x}) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, 1),$$

$$\textbf{Setting 2} : y = H_1(\boldsymbol{\beta}_*^\top \mathbf{x}) + \frac{1}{\sqrt{2}}H_2(\boldsymbol{\beta}_*^\top \mathbf{x}).$$

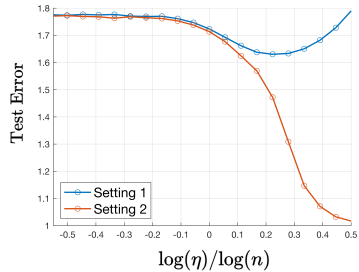
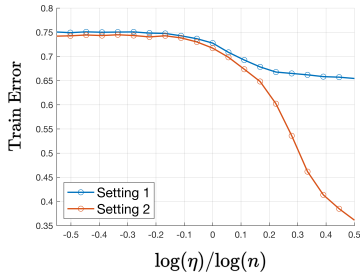


Simulation Results

We consider

$$\text{Setting 1 : } y = H_1(\beta_\star^\top \mathbf{x}) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, 1),$$

$$\text{Setting 2 : } y = H_1(\beta_\star^\top \mathbf{x}) + \frac{1}{\sqrt{2}}H_2(\beta_\star^\top \mathbf{x}).$$



Conclusion



Conclusion

- One-step gradient descent with step size $\eta \asymp n^\alpha$ can lead to feature learning.



Conclusion

- One-step gradient descent with step size $\eta \asymp n^\alpha$ can lead to feature learning.
- Learned features depend on the range of α .



Conclusion

- One-step gradient descent with step size $\eta \asymp n^\alpha$ can lead to feature learning.
- Learned features depend on the range of α .
- Unlike random features model, after one gradient update, the model can learn higher order polynomial components of the teacher function.

References



References I

- Ba, J., Erdogdu, M. A., Suzuki, T., Wang, Z., Wu, D., and Yang, G. (2022). High-dimensional asymptotics of feature learning: How one gradient step improves the representation. In *Advances in Neural Information Processing Systems*.
- Bombari, S., Kiyani, S., and Mondelli, M. (2023). Beyond the universal law of robustness: Sharper laws for random features and neural tangent kernels. In *International Conference on Machine Learning*.
- Bombari, S. and Mondelli, M. (2023). Stability, generalization and privacy: Precise analysis for random and NTK features. *arXiv preprint arXiv:2305.12100*.
- Clarté, L., Loureiro, B., Krzakala, F., and Zdeborová, L. (2023). On double-descent in uncertainty quantification in overparametrized models. In *International Conference on Artificial Intelligence and Statistics*.
- Dandi, Y., Krzakala, F., Loureiro, B., Pesce, L., and Stephan, L. (2023). Learning two-layer neural networks, one (giant) step at a time. *arXiv preprint arXiv:2305.18270*.
- Hassani, H. and Javanmard, A. (2022). The curse of overparametrization in adversarial training: Precise analysis of robust generalization for random features regression. *arXiv preprint arXiv:2201.05149*.
- Hu, H. and Lu, Y. M. (2023). Universality laws for high-dimensional learning with random features. *IEEE Transactions on Information Theory*, 69(3).
- Lee, D., Moniri, B., Huang, X., Dobriban, E., and Hassani, H. (2023). Demystifying disagreement-on-the-line in high dimensions. In *International Conference on Machine Learning*.
- Lin, L. and Dobriban, E. (2021). What causes the test error? going beyond bias-variance via ANOVA. *Journal of Machine Learning Research*, 22:155–1.
- Mei, S. and Montanari, A. (2022). The generalization error of random features regression: Precise asymptotics and the double descent curve. *Communications on Pure and Applied Mathematics*, 75(4):667–766.
- Moniri, B., Lee, D., Hassani, H., and Dobriban, E. (2023). A theory of non-linear feature learning with one gradient step in two-layer neural networks. *arXiv preprint arXiv:2310.07891*.
- Rahimi, A. and Recht, B. (2007). Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems*.

Questions?