

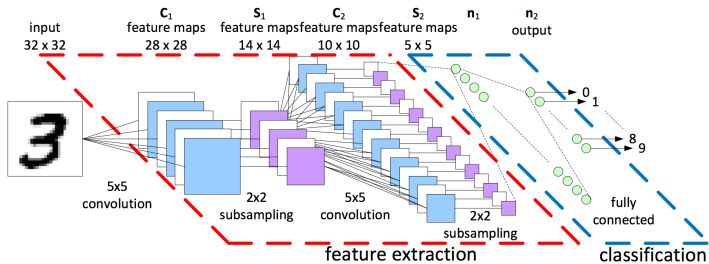
# Linearized two-layers neural network in high dimensions

Song Mei

Stanford University

May 26, 2019

Joint work with Andrea Montanari, Theodor Misiakiewicz, Behrooz Ghorbani



$$R(\Theta) = \min_{\Theta} \mathbb{E}[\ell(y, \mathbf{W}_1 \sigma \circ \mathbf{W}_2 \circ \sigma \circ \cdots \circ \mathbf{W}_k \circ \mathbf{x})].$$

### Empirical surprise of neural network [Zhang *et al.*, 2016]

- ▶ Over-parameterized regime.
- ▶ Optimization surprise: efficiently fit all the data.
- ▶ Generalization surprise: generalize well.

# Two-layers neural network

$$\hat{f}_N(\mathbf{x}; \Theta) = \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle), \quad \Theta = (a_1, \mathbf{w}_1, \dots, a_N, \mathbf{w}_N).$$

- ▶ Feature  $\mathbf{x} \in \mathbb{R}^d$ .
- ▶ Bottom layer weights  $\mathbf{w}_i \in \mathbb{R}^d$ ,  $i = 1, 2, \dots, N$ .
- ▶ Top layer weights  $a_i \in \mathbb{R}$ ,  $i = 1, 2, \dots, N$ .
- ▶ Over-parametrization:  $N$  large.

# Two-layers neural network

$$\hat{f}_N(\mathbf{x}; \Theta) = \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle), \quad \Theta = (a_1, \mathbf{w}_1, \dots, a_N, \mathbf{w}_N).$$

- ▶ Feature  $\mathbf{x} \in \mathbb{R}^d$ .
- ▶ Bottom layer weights  $\mathbf{w}_i \in \mathbb{R}^d$ ,  $i = 1, 2, \dots, N$ .
- ▶ Top layer weights  $a_i \in \mathbb{R}$ ,  $i = 1, 2, \dots, N$ .
- ▶ Over-parametrization:  $N$  large.

## Two-layers neural network

$$\hat{f}_N(\mathbf{x}; \Theta) = \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle), \quad \Theta = (a_1, \mathbf{w}_1, \dots, a_N, \mathbf{w}_N).$$

- ▶ Feature  $\mathbf{x} \in \mathbb{R}^d$ .
- ▶ Bottom layer weights  $\mathbf{w}_i \in \mathbb{R}^d$ ,  $i = 1, 2, \dots, N$ .
- ▶ Top layer weights  $a_i \in \mathbb{R}$ ,  $i = 1, 2, \dots, N$ .
- ▶ Over-parametrization:  $N$  large.

## Two-layers neural network

$$\hat{f}_N(\mathbf{x}; \Theta) = \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle), \quad \Theta = (a_1, \mathbf{w}_1, \dots, a_N, \mathbf{w}_N).$$

- ▶ Feature  $\mathbf{x} \in \mathbb{R}^d$ .
- ▶ Bottom layer weights  $\mathbf{w}_i \in \mathbb{R}^d$ ,  $i = 1, 2, \dots, N$ .
- ▶ Top layer weights  $a_i \in \mathbb{R}$ ,  $i = 1, 2, \dots, N$ .
- ▶ Over-parametrization:  $N$  large.

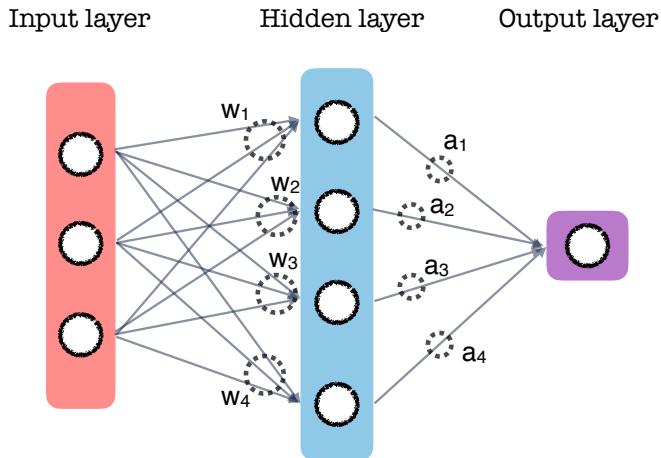
## Two-layers neural network

$$\hat{f}_N(\mathbf{x}; \Theta) = \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle), \quad \Theta = (a_1, \mathbf{w}_1, \dots, a_N, \mathbf{w}_N).$$

- ▶ Feature  $\mathbf{x} \in \mathbb{R}^d$ .
- ▶ Bottom layer weights  $\mathbf{w}_i \in \mathbb{R}^d$ ,  $i = 1, 2, \dots, N$ .
- ▶ Top layer weights  $a_i \in \mathbb{R}$ ,  $i = 1, 2, \dots, N$ .
- ▶ Over-parametrization:  $N$  large.



# Two-layers neural network



# Gradient flow with random initialization

Empirical risk: ( $n$ : # data;  $N$ : # neuron)

$$R_{n,N}(\Theta) = \hat{\mathbb{E}}_{\mathbf{x},n}[(\mathbf{y} - \hat{f}_N(\mathbf{x}; \Theta))^2]$$

Gradient flow, on empirical risk, with random initialization:

$$\begin{aligned}\dot{\Theta}(t) &= -\nabla R_{n,N}(\Theta(t)), \\ (a_i(0), \mathbf{w}_i(0)) &\sim_{i.i.d.} \mathbb{P}_{a,w}.\end{aligned}$$

# Convergence guarantees

Lemma (Global min. Not surprise. )

For  $N > n$ , we have

$$\inf_{\Theta} R_{n,N}(\Theta) = 0.$$

*There are many global minimizers with empirical risk 0.*

# Convergence guarantees

Lemma (Global min. Not surprise. )

For  $N > n$ , we have

$$\inf_{\Theta} R_{n,N}(\Theta) = 0.$$

*There are many global minimizers with empirical risk 0.*

But there are also local minimizers with non-zero risk.

## Convergence guarantees

Lemma (Global min. Not surprise. )

For  $N > n$ , we have

$$\inf_{\Theta} R_{n,N}(\Theta) = 0.$$

There are many global minimizers with empirical risk 0.

But there are also local minimizers with non-zero risk.

Theorem (The optimization surprise. )

For  $N \gg n^{1+c}$ , we have

$$\lim_{t \rightarrow \infty} R_{n,N}(\Theta(t)) = 0,$$

i.e., training loss converges to 0.

## Convergence guarantees

Lemma (Global min. Not surprise. )

For  $N > n$ , we have

$$\inf_{\Theta} R_{n,N}(\Theta) = 0.$$

*There are many global minimizers with empirical risk 0.*

But there are also local minimizers with non-zero risk.

Theorem (The optimization surprise. )

For  $N \gg n^{1+c}$ , we have

$$\lim_{t \rightarrow \infty} R_{n,N}(\Theta(t)) = 0,$$

*i.e., training loss converges to 0.*

Under what assumptions?

## Three variants of the convergence theorem

Gradient flow ( $n$ : # data;  $N$ : # neuron):

$$\begin{aligned}\dot{\Theta}(t) &= -\nabla \hat{\mathbb{E}}_{x,n}[(y - \hat{f}_N(x; \Theta(t)))^2], \\ \mathbf{w}_i(0) &\sim_{i.i.d.} \mathcal{N}(\mathbf{0}, \mathbf{I}_d/d).\end{aligned}$$

Theorem: for  $N$  large enough, we have

$$\lim_{t \rightarrow \infty} R_{n,N}(\Theta(t)) = 0.$$

### Random feature (RF) regime

$$\hat{f}_N(x; \Theta) = \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, x \rangle), \quad a_i(0) \sim_{i.i.d.} \mathcal{N}(0, 1/N^2).$$

[Andoni *et al.*, 2014], [Danialy, 2017], [Yehudai and Shamir, 2019] ...

# Three variants of the convergence theorem

Gradient flow ( $n$ : # data;  $N$ : # neuron):

$$\begin{aligned}\dot{\Theta}(t) &= -\nabla \hat{\mathbb{E}}_{\mathbf{x}, n}[(\mathbf{y} - \hat{f}_N(\mathbf{x}; \Theta(t)))^2], \\ \mathbf{w}_i(0) &\sim_{i.i.d.} \mathcal{N}(\mathbf{0}, \mathbf{I}_d/d).\end{aligned}$$

Theorem: for  $N$  large enough, we have

$$\lim_{t \rightarrow \infty} R_{n, N}(\Theta(t)) = 0.$$

## Neural tangent (NT) regime

$$\hat{f}_N(\mathbf{x}; \Theta) = \frac{1}{\sqrt{N}} \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle), \quad a_i(0) \sim_{i.i.d.} \mathcal{N}(0, 1).$$

[Jacot *et al.*, 2018], [Du *et al.*, 2018], [Du *et al.*, 2018], [Allen-Zhu *et al.*, 2018], [Zou *et al.*, 2018]...



## Three variants of the convergence theorem

Gradient flow ( $n$ : # data;  $N$ : # neuron):

$$\begin{aligned}\dot{\Theta}(t) &= -\nabla \hat{\mathbb{E}}_{\mathbf{x}, n}[(\mathbf{y} - \hat{f}_N(\mathbf{x}; \Theta(t)))^2], \\ \mathbf{w}_i(0) &\sim_{i.i.d.} \mathcal{N}(\mathbf{0}, \mathbf{I}_d/d).\end{aligned}$$

Theorem: for  $N$  large enough, we have \*

$$\lim_{t \rightarrow \infty} R_{n, N}(\Theta(t)) = 0.$$

### Mean field (MF) regime

$$\hat{f}_N(\mathbf{x}; \Theta) = \frac{1}{N} \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle), \quad a_i(0) \sim_{i.i.d.} \mathcal{N}(0, 1).$$

[Mei *et al.*, 2018], [Rotskoff and Vanden-Eijden, 2018], [Chizat and Bach, 2018]...

# Three variants of the convergence theorem

## Random feature (RF) regime

$$\hat{f}_N(\mathbf{x}; \Theta) = \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle), \quad a_i \sim \mathcal{N}(0, 1/N^2).$$

## Neural tangent (NT) regime

$$\hat{f}_N(\mathbf{x}; \Theta) = \frac{1}{\sqrt{N}} \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle), \quad a_i \sim \mathcal{N}(0, 1).$$

## Mean field (MF) regime

$$\hat{f}_N(\mathbf{x}; \Theta) = \frac{1}{N} \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle), \quad a_i \sim \mathcal{N}(0, 1).$$

... but different behavior of dynamics

## Random feature (RF) regime

$$\hat{f}_N(\mathbf{x}; \Theta) = \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle), \quad a_i \sim \mathcal{N}(0, 1/N^2).$$

- ▶ The limiting dynamics is linear (effectively only  $\mathbf{a}$  is updated).
- ▶ Prediction function: kernel ridge regression with kernel

$$k_{\text{RF}}(\mathbf{x}, \mathbf{z}) = \hat{\mathbb{E}}_{\mathbf{w}, N}[\sigma(\langle \mathbf{w}, \mathbf{z} \rangle)\sigma(\langle \mathbf{w}, \mathbf{x} \rangle)].$$

[Andoni *et al.*, 2014], [Danialy, 2017], [Yehudai and Shamir, 2019] ...

... but different behavior of dynamics

## Neural tangent (NT) regime

$$\hat{f}_N(\mathbf{x}; \Theta) = \frac{1}{\sqrt{N}} \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle), \quad a_i \sim \mathcal{N}(0, 1).$$

- ▶ The limiting dynamics is linear (the change of  $\Theta$  is small).
- ▶ Prediction function: kernel ridge regression with kernel

$$k_{\text{NT}}(\mathbf{x}, \mathbf{z}) = \hat{\mathbb{E}}_{\mathbf{w}, N}[\sigma'(\langle \mathbf{w}, \mathbf{x} \rangle) \sigma'(\langle \mathbf{w}, \mathbf{z} \rangle)] \langle \mathbf{x}, \mathbf{z} \rangle + k_{\text{RF}}(\mathbf{x}, \mathbf{z}).$$

[Jacot *et al.*, 2018], [Du *et al.*, 2018], [Du *et al.*, 2018], [Allen-Zhu *et al.*, 2018], [Zou *et al.*, 2018]...

... but different behavior of dynamics

## Mean field (MF) regime

$$\hat{f}_N(\mathbf{x}; \Theta) = \frac{1}{N} \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle), \quad a_i \sim \mathcal{N}(0, 1).$$

- ▶ The limiting dynamics is non-linear (both  $\mathbf{a}$  and  $\mathbf{W}$  are updated).
- ▶ Distributional dynamics:  
$$\partial_t \rho_t(\mathbf{a}, \mathbf{w}) = \nabla \cdot (\rho \nabla \Psi(\mathbf{a}, \mathbf{w}; \rho_t)) + \beta^{-1} \Delta \rho_t.$$
- ▶ Prediction function:  $\hat{f}(\mathbf{x}; \rho_\infty) = \int \mathbf{a} \sigma(\langle \mathbf{w}, \mathbf{x} \rangle) \rho_\infty(\mathrm{d}\mathbf{a} \mathrm{d}\mathbf{w})$ .

[Mei *et al.*, 2018], [Rotskoff and Vanden-Eijnden, 2018], [Chizat and Bach, 2018], [Sirignano and Spiliopoulos, 2018]...

Optimization: 0 training loss.

Test risk = training loss + generalization risk.

Today: generalization.

Optimization: 0 training loss.

Test risk = training loss + generalization risk.

Today: generalization.

# Generalization theory for kernel methods

- ▶ Traditional theory: assume  $f_\star \in \text{RKHS}$ , then kernel ridge regression generalize well.
- ▶ Problem: in high dimension, RKHS is a very small space.

Today: in high dimension, kernel methods (RF and NT) don't generalize well.



# Generalization theory for kernel methods

- ▶ Traditional theory: assume  $f_\star \in \text{RKHS}$ , then kernel ridge regression generalize well.
- ▶ Problem: in high dimension, RKHS is a very small space.

Today: in high dimension, kernel methods (RF and NT) don't generalize well.

# Generalization theory for kernel methods

- ▶ Traditional theory: assume  $f_* \in \text{RKHS}$ , then kernel ridge regression generalize well.
- ▶ Problem: in high dimension, RKHS is a very small space.

Today: in high dimension, kernel methods (RF and NT) don't generalize well.

## Setting 1: $N$ finite, $n$ infinite

Distribution:

$$\mathbf{x} \in \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})), \quad \mathbf{y} = f_*(\mathbf{x}), \quad f_* \in L^2(\mathbb{S}^d(\sqrt{d})).$$

Two classes of linearized neural network: ( $\mathbf{w}_i \sim \text{Unif}(\mathbb{S}^d)$ )

$$\mathcal{F}_{\text{RF},N}(\mathbf{W}) = \left\{ f = \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) : a_i \in \mathbb{R}, i \in [N] \right\},$$

$$\mathcal{F}_{\text{NT},N}(\mathbf{W}) = \left\{ f = \sum_{i=1}^N \sigma'(\langle \mathbf{w}_i, \mathbf{x} \rangle) \langle \mathbf{a}_i, \mathbf{x} \rangle : \mathbf{a}_i \in \mathbb{R}^d, i \in [N] \right\}.$$

Mild assumptions on  $\sigma$  (universal approximation, growth not too fast).

Lower bound:  $N$  finite,  $n$  infinite

$$\mathcal{F}_{\text{RF},N}(\mathbf{W}) = \left\{ f = \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) : a_i \in \mathbb{R}, i \in [N] \right\}.$$

Theorem (Ghorbani, Mei, Misiakiwics, Montanari, 2019)

Assume  $N = O_d(d^{\ell-\delta})$ , and  $(\mathbf{w}_i)_{i \in [N]} \sim \text{Unif}(\mathbb{S}^d)$ , we have

$$\inf_{f \in \mathcal{F}_{\text{RF},N}(\mathbf{W})} \mathbb{E}_{\mathbf{x}}[(f_{\star}(\mathbf{x}) - f(\mathbf{x}))^2] \geq \|P_{>\ell} f_{\star}\|_{L^2}^2 + o_{d,\mathbb{P}}(\|f_{\star}\|_2^2),$$

where  $P_{>\ell}$  is the projection operator orthogonal to the space of degree- $\ell$  polynomials.

Example: for  $f_{\star}(x) = x_1^2 - 1$ , we have  $P_{>2} f_{\star} \approx f_{\star}$ . Then random feature regression with  $N = O_d(d^{2-\delta})$  neuron achieves trivial risk, which is  $\|f_{\star}\|_{L^2}^2$ .

Lower bound:  $N$  finite,  $n$  infinite

$$\mathcal{F}_{\text{RF},N}(\mathbf{W}) = \left\{ f = \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) : a_i \in \mathbb{R}, i \in [N] \right\}.$$

Theorem (Ghorbani, Mei, Misiakiwics, Montanari, 2019)

Assume  $N = O_d(d^{\ell-\delta})$ , and  $(\mathbf{w}_i)_{i \in [N]} \sim \text{Unif}(\mathbb{S}^d)$ , we have

$$\inf_{f \in \mathcal{F}_{\text{RF},N}(\mathbf{W})} \mathbb{E}_{\mathbf{x}}[(f_{\star}(\mathbf{x}) - f(\mathbf{x}))^2] \geq \|P_{>\ell} f_{\star}\|_{L^2}^2 + o_{d,\mathbb{P}}(\|f_{\star}\|_2^2),$$

where  $P_{>\ell}$  is the projection operator orthogonal to the space of degree- $\ell$  polynomials.

Example: for  $f_{\star}(\mathbf{x}) = x_1^2 - 1$ , we have  $P_{>2} f_{\star} \approx f_{\star}$ . Then random feature regression with  $N = O_d(d^{2-\delta})$  neuron achieves trivial risk, which is  $\|f_{\star}\|_{L^2}^2$ .

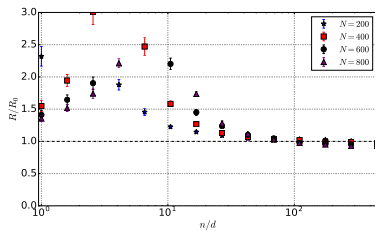
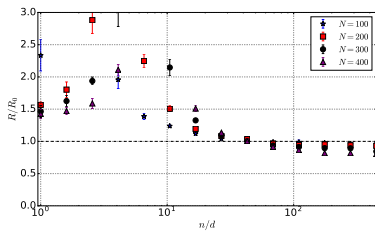


Figure: Test risk for learning  $f(\mathbf{x}) = x_1^2 - 1$ ,  $d = 50$  and  $d = 100$ .

## Similar result for NT

$$\mathcal{F}_{\text{NT},N}(\mathbf{W}) = \left\{ f = \sum_{i=1}^N \sigma'(\langle \mathbf{w}_i, \mathbf{x} \rangle) \langle \mathbf{a}_i, \mathbf{x} \rangle : \mathbf{a}_i \in \mathbb{R}^d, i \in [N] \right\}.$$

Theorem (Ghorbani, Mei, Misiakiwics, Montanari, 2019)

Assume  $N = O_d(d^{\ell-\delta})$ , and  $(\mathbf{w}_i)_{i \in [N]} \sim \text{Unif}(\mathbb{S}^d)$ , we have

$$\inf_{f \in \mathcal{F}_{\text{RF},N}(\mathbf{W})} \mathbb{E}_{\mathbf{x}}[(f_{\star}(\mathbf{x}) - f(\mathbf{x}))^2] \geq \|P_{>\ell+1} f_{\star}\|_{L^2}^2 + o_{d,\mathbb{P}}(\|f_{\star}\|_2^2),$$

where  $P_{>\ell+1}$  is the projection operator orthogonal to the space of degree- $(\ell + 1)$  polynomials.

Example: for  $f_{\star}(x) = x_1^3 - x_1$ , we have  $P_{>3} f_{\star} \approx f_{\star}$ . Then random feature regression with  $N = O_d(d^{2-\delta})$  neuron achieves trivial risk, which is  $\|f_{\star}\|_{L^2}^2$ .

## Similar result for NT

$$\mathcal{F}_{\text{NT},N}(\mathbf{W}) = \left\{ f = \sum_{i=1}^N \sigma'(\langle \mathbf{w}_i, \mathbf{x} \rangle) \langle \mathbf{a}_i, \mathbf{x} \rangle : \mathbf{a}_i \in \mathbb{R}^d, i \in [N] \right\}.$$

Theorem (Ghorbani, Mei, Misiakiwics, Montanari, 2019)

Assume  $N = O_d(d^{\ell-\delta})$ , and  $(\mathbf{w}_i)_{i \in [N]} \sim \text{Unif}(\mathbb{S}^d)$ , we have

$$\inf_{f \in \mathcal{F}_{\text{RF},N}(\mathbf{W})} \mathbb{E}_{\mathbf{x}}[(f_{\star}(\mathbf{x}) - f(\mathbf{x}))^2] \geq \|P_{>\ell+1} f_{\star}\|_{L^2}^2 + o_{d,\mathbb{P}}(\|f_{\star}\|_2^2),$$

where  $P_{>\ell+1}$  is the projection operator orthogonal to the space of degree- $(\ell + 1)$  polynomials.

Example: for  $f_{\star}(\mathbf{x}) = x_1^3 - x_1$ , we have  $P_{>3} f_{\star} \approx f_{\star}$ . Then random feature regression with  $N = O_d(d^{2-\delta})$  neuron achieves trivial risk, which is  $\|f_{\star}\|_{L^2}^2$ .



## Setting 2: $N$ infinite, $n$ finite

Distribution:

$$\mathbf{x}_i \in \text{Unif}(\mathbb{S}^{d-1}), \quad y_i = f_\star(\mathbf{x}_i), \quad f_\star \in L^2(\mathbb{S}^d(\sqrt{d})).$$

Predicting using regularized kernel ridge regression:

$$\hat{f}_\lambda(\mathbf{x}) = k(\mathbf{x}, \mathcal{X})(k(\mathcal{X}, \mathcal{X}) + \lambda \mathbf{I})^{-1} f_\star(\mathbf{x}),$$

where

$$k(\mathbf{x}_i, \mathbf{x}_j) = \mathbb{E}_{\mathbf{w} \sim \mathbb{S}^{d-1}} [\sigma(\langle \mathbf{w}, \mathbf{x}_i \rangle) \sigma(\langle \mathbf{w}, \mathbf{x}_j \rangle)].$$

Lower bound:  $N$  infinite,  $n$  finite

Theorem (Ghorbani, Mei, Misiakiwics, Montanari, 2019)

Assume  $n = O_d(d^{\ell-\delta})$ , we have

$$\inf_{\lambda} \mathbb{E}_{\mathbf{x}}[(f_{\star}(\mathbf{x}) - \hat{f}_{\lambda}(\mathbf{x}))^2] \geq \|P_{>\ell} f_{\star}\|_{L^2}^2 + o_d(\|f_{\star}\|_2^2),$$

where  $P_{>\ell}$  is the projection operator orthogonal to the space of degree- $\ell$  polynomials.

# Intuition behind these results

In high dimension, the correlation between a degree- $k$  Hermite polynomial and a random feature is very small

$$\mathbb{E}_{\mathbf{w}}[\text{He}_k(\mathbf{x}_1)\sigma(\langle \mathbf{w}, \mathbf{x} \rangle)] = O_d(1/d^k).$$

Also observed in [Danialy, 2016], [Bach, 2017].

# Implications & Conclusions

- ▶ In high dimension, even for simple function  $f(x) = x_1^k$ , it takes  $n, N = O_d(d^k)$  to learn it well using **linearized** neural network (kernel methods);
- ▶ ... while a neural network can learn it (conjecture to be efficiently) using  $n, N = O_d(1)$ .
- ▶ Neural network is more powerful than kernel methods.
- ▶ Future work: what class of functions neural network can learn efficiently.

# Implications & Conclusions

- ▶ In high dimension, even for simple function  $f(x) = x_1^k$ , it takes  $n, N = O_d(d^k)$  to learn it well using **linearized** neural network (kernel methods);
- ▶ ... while a neural network can learn it (conjecture to be efficiently) using  $n, N = O_d(1)$ .
- ▶ Neural network is more powerful than kernel methods.
- ▶ Future work: what class of functions neural network can learn efficiently.

# Implications & Conclusions

- ▶ In high dimension, even for simple function  $f(x) = x_1^k$ , it takes  $n, N = O_d(d^k)$  to learn it well using **linearized** neural network (kernel methods);
- ▶ ... while a neural network can learn it (conjecture to be efficiently) using  $n, N = O_d(1)$ .
- ▶ Neural network is more powerful than kernel methods.
- ▶ Future work: what class of functions neural network can learn efficiently.

# Implications & Conclusions

- ▶ In high dimension, even for simple function  $f(x) = x_1^k$ , it takes  $n, N = O_d(d^k)$  to learn it well using **linearized** neural network (kernel methods);
- ▶ ... while a neural network can learn it (conjecture to be efficiently) using  $n, N = O_d(1)$ .
- ▶ Neural network is more powerful than kernel methods.
- ▶ Future work: what class of functions neural network can learn efficiently.