# Thermodynamic limits for neural networks

Andrea Montanari

Stanford University

November 21, 2019



- $\blacktriangleright \ \{(y_i, x_i)\}_{i \leq n}$
- $x_i \sim {\sf Unif}(\mathbb{S}^{d-1}(\sqrt{d})) ext{ or } x_i \sim {\sf N}(0, {I}_d) ext{ , } d \gg 1$
- Response (regression)

$$y_i = f_*(x_i) + arepsilon_i\,, \qquad arepsilon_i \sim {\sf N}(0, au^2)$$

## Evaluating a model

▶ Loss function  $\ell : \mathbb{R} \times \mathbb{R} \to \mathbb{R}_{\geq 0}$  (e.g.  $\ell(y_1, y_2) = (y_1 - y_2)^2$ )

► Test error

$$R(f) = \mathbb{E}_{ ext{new}} \Big\{ \ \ell(y^{ ext{new}}, f(x^{ ext{new}})) \Big\}$$

## Evaluating a model

▶ Loss function  $\ell : \mathbb{R} \times \mathbb{R} \to \mathbb{R}_{\geq 0}$  (e.g.  $\ell(y_1, y_2) = (y_1 - y_2)^2$ )

► Test error

$$R(f) = \mathbb{E}_{ ext{new}} \Big\{ \ oldsymbol{\ell}(y^{ ext{new}}, f(x^{ ext{new}})) \Big\}$$

# Two-layers (fully-connected) neural networks

$$\mathcal{F}_{\mathsf{NN}}^N \equiv \left\{ f_{oldsymbol{ heta}}(oldsymbol{x}) = \sum_{i=1}^N a_i \, \sigma(\langle oldsymbol{w}_i, oldsymbol{x} 
angle) \, : \ oldsymbol{ heta} = (a_i, oldsymbol{w}_i)_{i \leq N}, \ a_i \in \mathbb{R}, oldsymbol{w}_i \in \mathbb{R}^d 
ight\}$$

#### • $\hat{\theta}_n$ computed on training sample

► Train/test error

$$egin{aligned} \widehat{R}_n(m{ heta}) &= rac{1}{n}\sum_{i=1}^n \ell(y_i, f_{m{ heta}}(x))\,, \ R(\hat{ heta}_n) &= \mathbb{E}_{ ext{new}}\{\ell(y^{ ext{new}}, f_{m{ heta}_n}(x^{ ext{new}}))\} \end{aligned}$$

# Two-layers (fully-connected) neural networks

$$\mathcal{F}_{\mathsf{NN}}^N \equiv \left\{ f_{oldsymbol{ heta}}(x) = \sum_{i=1}^N a_i \, \sigma(\langle oldsymbol{w}_i, x 
angle) : \hspace{0.2cm} oldsymbol{ heta} = (a_i, oldsymbol{w}_i)_{i \leq N}, \hspace{0.2cm} a_i \in \mathbb{R}, oldsymbol{w}_i \in \mathbb{R}^d 
ight\}$$

### • $\hat{\theta}_n$ computed on training sample

► Train/test error

$$egin{aligned} \widehat{R}_n(oldsymbol{ heta}) &= rac{1}{n}\sum_{i=1}^n \ell(y_i, f_{oldsymbol{ heta}}(x))\,, \ R(\hat{oldsymbol{ heta}}_n) &= \mathbb{E}_{ ext{new}}\{\ell(y^{ ext{new}}, f_{oldsymbol{ heta}_n}(x^{ ext{new}}))\} \end{aligned}$$

#### ▶ Classical limit in Stats/ML: $n \to \infty$ (large sample)

• Here:  $N \to \infty$  (large network)

#### ▶ Classical limit in Stats/ML: $n \to \infty$ (large sample)

▶ Here:  $N \to \infty$  (large network)

▶ Classical limit in Stats/ML:  $n \to \infty$  (large sample)

• Here:  $N \to \infty$  (large network)

▶ Classical limit in Stats/ML:  $n \to \infty$  (large sample)

• Here:  $N \to \infty$  (large network)

$$f_{oldsymbol{ heta}}(oldsymbol{x}) = \sum_{i=1}^N a_i\,\sigma(\langleoldsymbol{w}_i,oldsymbol{x}
angle)$$

$$f_{m{ heta}}(m{x}) = \int a \sigma(\langle m{w},m{x}
angle) \ 
ho( ext{d} a, ext{d} w) \,, \quad 
ho = rac{1}{N} \sum_{i=1}^N \delta_{Na_i,w_i}$$

$$egin{split} \mathcal{F}_{\mathsf{NN}}^{N} \subseteq \mathcal{F}_{\mathsf{NN}}^{N+1} \subseteq \cdots \subseteq \mathcal{F}_{\mathsf{NN}}^{\infty}\,, \ \mathcal{F}_{\mathsf{NN}}^{\infty} \equiv & \left\{f(x;
ho) = \int a\sigma(\langle w,x
angle)\,\,
ho(\mathsf{d} a,\mathsf{d} w) \quad 
ho \in \mathscr{P}(\mathbb{R}^{d+1})
ight\}. \end{split}$$

$$f_{oldsymbol{ heta}}({m x}) = \sum_{i=1}^N a_i\,\sigma(\langle {m w}_i,{m x}
angle)$$

$$f_{m{ heta}}(x) = \int a \sigma(\langle w, x 
angle) \ 
ho( ext{d} a, ext{d} w) \,, \quad 
ho = rac{1}{N} \sum_{i=1}^N \delta_{Na_i,w_i}$$

$$egin{split} \mathcal{F}_{\mathsf{NN}}^{N} &\subseteq \mathcal{F}_{\mathsf{NN}}^{N+1} \subseteq \cdots \subseteq \mathcal{F}_{\mathsf{NN}}^{\infty}\,, \ \mathcal{F}_{\mathsf{NN}}^{\infty} &\equiv \left\{f(x;
ho) = \int a\sigma(\langle w,x
angle) \;
ho(\mathsf{d} a,\mathsf{d} w) \quad 
ho \in \mathscr{P}(\mathbb{R}^{d+1})
ight\}. \end{split}$$

$$f_{oldsymbol{ heta}}(x) = \sum_{i=1}^N a_i\,\sigma(\langle oldsymbol{w}_i,x
angle)$$

$$f_{oldsymbol{ heta}}(x) = \int a \sigma(\langle w, x 
angle) \ 
ho(\mathrm{d} a, \mathrm{d} w) \,, \quad 
ho = rac{1}{N} \sum_{i=1}^N \delta_{Na_i, w_i} \,.$$

$$egin{split} \mathcal{F}_{\mathsf{NN}}^{N} \subseteq \mathcal{F}_{\mathsf{NN}}^{N+1} \subseteq \cdots \subseteq \mathcal{F}_{\mathsf{NN}}^{\infty}\,, \ \mathcal{F}_{\mathsf{NN}}^{\infty} \equiv \left\{f(x;
ho) = \int a\sigma(\langle w,x
angle)\,\,
ho(\mathrm{d} a,\mathrm{d} w) \quad 
ho \in \mathscr{P}(\mathbb{R}^{d+1})
ight\}. \end{split}$$

$$f_{oldsymbol{ heta}}(x) = \sum_{i=1}^N a_i\,\sigma(\langle oldsymbol{w}_i,x
angle)$$

$$f_{oldsymbol{ heta}}(x) = \int a \sigma(\langle w, x 
angle) \ 
ho(\mathrm{d} a, \mathrm{d} w) \,, \quad 
ho = rac{1}{N} \sum_{i=1}^N \delta_{Na_i, w_i} \,.$$

$$egin{split} \mathcal{F}_{\mathsf{NN}}^{N} \subseteq \mathcal{F}_{\mathsf{NN}}^{N+1} \subseteq \cdots \subseteq \mathcal{F}_{\mathsf{NN}}^{\infty}\,, \ \mathcal{F}_{\mathsf{NN}}^{\infty} \equiv \left\{f(x;
ho) = \int a\sigma(\langle w,x
angle)\,\,
ho(\mathsf{d} a,\mathsf{d} w) \quad 
ho \in \mathscr{P}(\mathbb{R}^{d+1})
ight\}. \end{split}$$

$$\widehat{R}_n(
ho) = rac{1}{n}\sum_{i=1}^n \ell(y_i,f(x_i;
ho))$$

- **Good news:** Convex
- Bad news: Infinite dimensional
- Question: Can we optimize it efficiently?
   e.g. approximate ρ by finitely many points

$$\widehat{R}_n(
ho) = rac{1}{n}\sum_{i=1}^n \ell(y_i,f(m{x}_i;
ho))$$

#### ► Good news: Convex

#### Bad news: Infinite dimensional

Question: Can we optimize it efficiently?
 e.g. approximate ρ by finitely many points

$$\widehat{R}_n(
ho) = rac{1}{n}\sum_{i=1}^n \ell(y_i,f(m{x}_i;
ho))$$

► Good news: Convex

#### Bad news: Infinite dimensional

Question: Can we optimize it efficiently?
 e.g. approximate ρ by finitely many points

$$\widehat{R}_n(
ho) = rac{1}{n}\sum_{i=1}^n \ell(y_i,f(x_i;
ho))$$

- Good news: Convex
- Bad news: Infinite dimensional
- Question: Can we optimize it efficiently?
   e.g. approximate ρ by finitely many points

$$\widehat{R}_n(
ho) = rac{1}{n}\sum_{i=1}^n \ell(y_i,f(x_i;
ho))$$

- Good news: Convex
- Bad news: Infinite dimensional
- Question: Can we optimize it efficiently?
   e.g. approximate ρ by finitely many points

Another approach: lazy limit

- ▶ Idea: Linearize around a random initialization
- $\blacktriangleright \Rightarrow \text{Linear class } \mathcal{F}_{\mathsf{NT}}$

## Neural tangent

$$egin{aligned} \mathcal{F}_{\mathsf{NT}}^N(oldsymbol{W}) &\equiv \left\{f(oldsymbol{x}) = \sum_{i=1}^N \langle oldsymbol{a}_i, oldsymbol{x} 
angle \sigma'(\langle oldsymbol{w}_i, oldsymbol{x} 
angle) \ : & oldsymbol{a}_i \in \mathbb{R}^d \ orall i \leq N 
ight\}, \ & oldsymbol{W} = [oldsymbol{w}_1, \dots, oldsymbol{w}_N] \quad oldsymbol{w}_i \sim_{iid} \mathsf{Unif}(\mathbb{S}^{d-1}(1)) \end{aligned}$$

# • 'Tangent' to $\mathcal{F}_{\mathsf{NN}}^N$ around a random point

• Good approximation: (i)  $N \gg n$ ; (ii) Parameters's scaling

[Jacot, Gabriel, Hongler, 2018; Du, Zhai, Poczos, Singh 2018; Allen-Zhu, Li, Song 2018; Chizat, Bach, 2019; Arora, Du, Hu, Li, Salakhutdinov, Wang, 2019; Oymak, Soltanolkotabi, 2019; ...]

## Neural tangent

$$egin{aligned} \mathcal{F}_{\mathsf{NT}}^N(oldsymbol{W}) &\equiv \left\{f(oldsymbol{x}) = \sum_{i=1}^N \langle oldsymbol{a}_i, oldsymbol{x} 
angle \sigma'(\langle oldsymbol{w}_i, oldsymbol{x} 
angle) \ : & oldsymbol{a}_i \in \mathbb{R}^d \ orall i \leq N 
ight\}, \ & oldsymbol{W} = [oldsymbol{w}_1, \dots, oldsymbol{w}_N] \quad oldsymbol{w}_i \sim_{iid} \mathsf{Unif}(\mathbb{S}^{d-1}(1)) \end{aligned}$$

- 'Tangent' to  $\mathcal{F}_{NN}^N$  around a random point
- Good approximation: (i)  $N \gg n$ ; (ii) Parameters's scaling

[Jacot, Gabriel, Hongler, 2018; Du, Zhai, Poczos, Singh 2018; Allen-Zhu, Li, Song 2018; Chizat, Bach, 2019; Arora, Du, Hu, Li, Salakhutdinov, Wang, 2019; Oymak, Soltanolkotabi, 2019; ...]

## A simpler linear model

#### Random features model

[Neal, 1996; Balcan, Blum, Vempala 2006; Rahimi, Recht; 2008; Bach, 2016; ...]



#### Which phenomena are captured linear models?

## Outline

1) What is captured by linear models?

2 What is not captured by linear models?



Ghorbani, Mei, Misiakiewicz, M, arXiv:1904.12191, 1906.08899 Mei, M, arXiv:1908.05355 M, Ruan, Sohn, Yan, arXiv:1911.01544

#### What is captured by linear models?

What is captured by linear models?

- Interpolation phase transition
- Double descent

## Using deep nets to fit noise



[Zhang, Bengio, Hardt, Recht, Vinyals, 2016]

# Interpolating, yet not overfitting



- ▶ MNIST (subset): 4,000 images in 10 different classes.
- ▶ 2-layers Neural Net. Square loss.

[Belkin, Hsu, Ma, Mandal, 2018]

# Interpolating, yet not overfitting



- ▶ MNIST: 50,000 images in 2 different classes.
- ▶ 5-layers Neural Net. Quadratic hinge loss.

[Spigler, Geiger, d'Ascoli, Sagun, Biroli, Wyart, 2018]

## The double-descent curve



[Belkin, Hsu, Ma, Mandal, 2018]



Peak at the interpolation threshold

- ✓ Global minimum in the overparametrized regime
- ✓ Monotone decreasing in the overparametrized regime
- ✓ Optimal regularization  $\rightarrow$  0.



#### Peak at the interpolation threshold

- ✓ Global minimum in the overparametrized regime
- ✓ Monotone decreasing in the overparametrized regime
- ✓ Optimal regularization  $\rightarrow$  0.



- $\checkmark$  Peak at the interpolation threshold
- ✓ Global minimum in the overparametrized regime
- ✓ Monotone decreasing in the overparametrized regime
   ✓ Optimal regularization → 0.



- $\checkmark$  Peak at the interpolation threshold
- ✓ Global minimum in the overparametrized regime
- ✓ Monotone decreasing in the overparametrized regime
   ✓ Optimal regularization → 0.



- $\checkmark$  Peak at the interpolation threshold
- ✓ Global minimum in the overparametrized regime
- $\checkmark$  Monotone decreasing in the overparametrized regime
- ✓ Optimal regularization  $\rightarrow$  0.
### Random features model

$$\mathcal{F}_{\mathsf{RF}}(\mathit{W}) \equiv \left\{ f(x) = \sum_{i=1}^N a_i \, \sigma(\langle w_i, x 
angle) \, : \quad a_i \in \mathbb{R} \; orall i \leq N 
ight\}.$$

• 
$$W = [w_1, \ldots, w_N]$$
 random;  $w_i \sim \text{Unif}(\mathbb{S}^{d-1}(1))$ , i.i.d.

# Ridge regression

$$\widehat{a}(\lambda) = \mathrm{argmin}_{oldsymbol{a} \in \mathbb{R}^N} \left\{ \widehat{\mathbb{E}}_n \Big[ \Big( y - \sum_{i=1}^N a_i \sigma(\langle oldsymbol{w}_i, oldsymbol{x} 
angle \Big)^2 \Big] + rac{N\lambda}{d} \, \|oldsymbol{a}\|_2^2 
ight\}$$



Random features ridge regression (as above.)

▶ Proportional regime  $n, N, d \rightarrow \infty$ ,

$$rac{N}{d} o \psi_1, \quad rac{n}{d} o \psi_2\,.$$

$$\blacktriangleright \ y_i = f_*(x_i) + \varepsilon_i$$

 $\blacktriangleright \ f_*(x) = \langle \boldsymbol{\beta}_0, x \rangle$ 

(Higher order terms  $\rightarrow$  See paper)

### Setting

Random features ridge regression (as above.)

▶ Proportional regime  $n, N, d \rightarrow \infty$ ,

$$rac{N}{d} o \psi_1, \quad rac{n}{d} o \psi_2$$
 .

$$\blacktriangleright \ y_i = f_*(x_i) + \varepsilon_i$$

 $\blacktriangleright \ f_*(x) = \langle \boldsymbol{\beta}_0, x \rangle$ 

(Higher order terms  $\rightarrow$  See paper)

### Setting

Random features ridge regression (as above.)

▶ Proportional regime  $n, N, d \rightarrow \infty$ ,

$$rac{N}{d} o \psi_1, \quad rac{n}{d} o \psi_2$$
 .

$$\begin{array}{ll} \blacktriangleright & y_i = f_*(x_i) + \varepsilon_i \\ \hline & f_*(x) = \langle \beta_0, x \rangle \end{array} \tag{Higher order terms} \rightarrow \text{See paper}$$



Random features ridge regression (as above.)

▶ Proportional regime  $n, N, d \rightarrow \infty$ ,

$$rac{N}{d} o \psi_1, \quad rac{n}{d} o \psi_2$$
 .

$$\begin{array}{ll} \blacktriangleright & y_i = f_*(x_i) + \varepsilon_i \\ \blacktriangleright & f_*(x) = \langle \beta_0, x \rangle \end{array} \tag{Higher order terms} \rightarrow \text{See paper}) \end{array}$$

### Precise asymptotics

#### Theorem (Mei, M. 2019)

Assume  $f_*(x) = \langle eta_0, x 
angle$  and define (for  $G \sim \mathsf{N}(0,1)$ )

$$b_*^2 = \mathbb{E}[\sigma(G)^2] - \mathbb{E}[\sigma(G)]^2 - b_1^2, \hspace{0.2cm} b_1 = \mathbb{E}[G\sigma(G)], \hspace{0.2cm} \zeta \equiv rac{b_1^2}{b_*^2}$$

Then, for any  $\lambda > 0$ , we have

 $R_{\mathsf{RF}}(\widehat{f}_{\lambda}) = \|oldsymbol{eta}_0\|_2^2 \mathscr{B}(\zeta,\psi_1,\psi_2,\lambda/\overline{b}_*^2) + au^2 \mathscr{V}(\zeta,\psi_1,\psi_2,\lambda/\overline{b}_*^2) + o_d(1),$ 

where  $\mathscr{B}(\zeta,\psi_1,\psi_2,\overline{\lambda}),\ \mathscr{V}(\zeta,\psi_1,\psi_2,\overline{\lambda})$  are explicitly given below.

Variance computed in [Hastie, M, Rosset, Tibshirani, 2019]

### Precise asymptotics

#### Theorem (Mei, M. 2019)

Assume  $f_*(x) = \langle eta_0, x 
angle$  and define (for  $G \sim \mathsf{N}(0,1)$ )

$$b_*^2 = \mathbb{E}[\sigma(G)^2] - \mathbb{E}[\sigma(G)]^2 - b_1^2, \hspace{0.2cm} b_1 = \mathbb{E}[G\sigma(G)], \hspace{0.2cm} \zeta \equiv rac{b_1^2}{b_x^2}$$

Then, for any  $\lambda > 0$ , we have

$$R_{\mathsf{RF}}(\widehat{f}_{\lambda}) = \|oldsymbol{eta}_0\|_2^2 \mathscr{B}(\zeta,\psi_1,\psi_2,\lambda/\overline{b}_*^2) + au^2 \mathscr{V}(\zeta,\psi_1,\psi_2,\lambda/\overline{b}_*^2) + o_d(1),$$

where  $\mathscr{B}(\zeta, \psi_1, \psi_2, \overline{\lambda})$ ,  $\mathscr{V}(\zeta, \psi_1, \psi_2, \overline{\lambda})$  are explicitly given below.

Variance computed in [Hastie, M, Rosset, Tibshirani, 2019]

### Explicit formulae

Let  $(\nu_1(\xi), \nu_2(\xi))$  be the unique solution of

$$\begin{split} \nu_1 &= \psi_1 \Big( -\xi - \nu_2 - \frac{\zeta^2 \nu_2}{1 - \zeta^2 \nu_1 \nu_2} \Big)^{-1} , \\ \nu_2 &= \psi_2 \Big( -\xi - \nu_1 - \frac{\zeta^2 \nu_1}{1 - \zeta^2 \nu_1 \nu_2} \Big)^{-1} ; \end{split}$$

Let

$$\chi \equiv 
u_1 ( \boldsymbol{i} ( \psi_1 \psi_2 \overline{\lambda} )^{1/2} ) \cdot 
u_2 ( \boldsymbol{i} ( \psi_1 \psi_2 \overline{\lambda} )^{1/2} ),$$

and

$$\begin{split} \mathscr{E}_{0}(\zeta,\psi_{1},\psi_{2},\overline{\lambda}) &\equiv -\chi^{5}\zeta^{6} + 3\chi^{4}\zeta^{4} + (\psi_{1}\psi_{2} - \psi_{2} - \psi_{1} + 1)\chi^{3}\zeta^{6} - 2\chi^{3}\zeta^{4} - 3\chi^{3}\zeta^{2} \\ &+ (\psi_{1} + \psi_{2} - 3\psi_{1}\psi_{2} + 1)\chi^{2}\zeta^{4} + 2\chi^{2}\zeta^{2} + \chi^{2} + 3\psi_{1}\psi_{2}\chi\zeta^{2} - \psi_{1}\psi_{2} , \\ \mathscr{E}_{1}(\zeta,\psi_{1},\psi_{2},\overline{\lambda}) &\equiv \psi_{2}\chi^{3}\zeta^{4} - \psi_{2}\chi^{2}\zeta^{2} + \psi_{1}\psi_{2}\chi\zeta^{2} - \psi_{1}\psi_{2} , \\ \mathscr{E}_{2}(\zeta,\psi_{1},\psi_{2},\overline{\lambda}) &\equiv \chi^{5}\zeta^{6} - 3\chi^{4}\zeta^{4} + (\psi_{1} - 1)\chi^{3}\zeta^{6} + 2\chi^{3}\zeta^{4} + 3\chi^{3}\zeta^{2} + (-\psi_{1} - 1)\chi^{2}\zeta^{4} - 2\chi^{2}\zeta^{2} - \chi^{2} . \end{split}$$

We then have

$$\mathscr{B}(\zeta,\psi_1,\psi_2,\overline{\lambda}) \equiv \ \frac{\mathscr{E}_1(\zeta,\psi_1,\psi_2,\overline{\lambda})}{\mathscr{E}_0(\zeta,\psi_1,\psi_2,\overline{\lambda})} , \qquad \mathscr{V}(\zeta,\psi_1,\psi_2,\overline{\lambda}) \equiv \ \frac{\mathscr{E}_2(\zeta,\psi_1,\psi_2,\overline{\lambda})}{\mathscr{E}_0(\zeta,\psi_1,\psi_2,\overline{\lambda})}$$

Random matrix theory for kernel inner product random matrices

[Cheng, Singer, 2013; Do, Vu 2013; Fan, M. 2019; Penington, Wohra, 2017;...]



Peak at the interpolation threshold

- ✓ Global minimum in the overparametrized regime
- ✓ Monotone decreasing in the overparametrized regime



- $\checkmark$  Peak at the interpolation threshold
- ✓ Global minimum in the overparametrized regime
- $\checkmark$  Monotone decreasing in the overparametrized regime

#### Insigths



- Singularity at the interpolation threshold
- Minimum risk at extreme overparametrization  $N/n \to \infty$ .



Same at  $\lambda > 0$  fixed: Minimum at  $N/n \to \infty$ .



- High SNR: Minimum at  $\lambda = 0+$ .
- Low SNR: Minimum at  $\lambda > 0$ .

#### Beyond square loss?

### **Binary classification**

▶ Data 
$$(y_i, x_i) \in \{+1, -1\} imes \mathbb{R}^d, \quad x_i \sim \mathsf{N}(0, I_d)$$
$$\mathbb{P}(y_i = +1 | x_i) = f_*(\langle \beta_0, x_i \rangle)$$

► Two-layer network

$$\hat{f}_a(x) = ext{sign} \left\{ \sum_{i=1}^N a_i \sigma(\langle w_i, x 
angle) 
ight\}$$

Prediction error

$$\mathrm{Err}_n \equiv \mathbb{P}(y^{ ext{new}} \widehat{f}_{\widehat{a}}(x^{ ext{new}}) \leq 0)$$
 ,

### **Binary classification**

▶ Data 
$$(y_i, x_i) \in \{+1, -1\} imes \mathbb{R}^d, \quad x_i \sim \mathsf{N}(0, I_d)$$
 $\mathbb{P}(y_i = +1 | x_i) = f_*(\langle eta_0, x_i \rangle)$ 

► Two-layer network

$$\hat{f}_{m{a}}(m{x}) = ext{sign} \left\{ \sum_{i=1}^{N} a_i \sigma(\langle m{w}_i, m{x} 
angle) 
ight\}$$

Prediction error

$$\mathrm{Err}_n \equiv \mathbb{P}(y^{\mathrm{new}} \widehat{f}_{\widehat{a}}(x^{\mathrm{new}}) \leq 0) \, ,$$

### **Binary classification**

▶ Data 
$$(y_i, x_i) \in \{+1, -1\} imes \mathbb{R}^d, \quad x_i \sim \mathsf{N}(0, I_d)$$
 $\mathbb{P}(y_i = +1 | x_i) = f_*(\langle eta_0, x_i \rangle)$ 

► Two-layer network

$$\hat{f}_{oldsymbol{a}}(oldsymbol{x}) = ext{sign} \left\{ \sum_{i=1}^N a_i \sigma(\langle oldsymbol{w}_i, oldsymbol{x} 
angle) 
ight\}$$

Prediction error

$$\mathrm{Err}_n \equiv \mathbb{P}(y^{\mathrm{new}} \widehat{f}_{\widehat{a}}(x^{\mathrm{new}}) \leq 0) \, ,$$

### Max-margin classification

#### Random features:

$$oldsymbol{W} = [oldsymbol{w}_1, \dots, oldsymbol{w}_N] \quad oldsymbol{w}_i \sim_{iid} \mathsf{Unif}(\mathbb{S}^{d-1}(1))$$

#### Maximize the margin:

$$egin{array}{cc} ext{maximize} & \min_{i\leq n} \, y_i \hat{f}_{m{a}}(x_i) = \min_{i\leq n} \, y_i \langle \, m{a}, \sigma(\,m{W}x_i) 
angle \, , \ & ext{subject to} & \| \, m{a} \|_2 \leq 1 \, . \end{array}$$

# $\max ig\{ \min_{i \leq n} \, y_i \langle a, \sigma({\it W} x_i) angle \, : \, \, \|a\| \leq 1 ig\}$

- Can't use random matrix theory
- $\sigma(Wx_i)$  has correlated entries
- ▶ Non-gaussian

$$\max ig\{ \min_{i \leq n} \ y_i \langle a, \sigma( oldsymbol{W} x_i) 
angle \ : \ \| oldsymbol{a} \| \leq 1 ig\}$$

#### Can't use random matrix theory

- $\sigma(Wx_i)$  has correlated entries
- ► Non-gaussian

$$\max ig\{ \min_{i \leq n} \ y_i \langle a, \sigma( oldsymbol{W} x_i) 
angle \ : \ \| oldsymbol{a} \| \leq 1 ig\}$$

- Can't use random matrix theory
- $\sigma(Wx_i)$  has correlated entries
- ▶ Non-gaussian

$$\maxig\{\min_{i\leq n}\ y_i\langle a,\sigma({oldsymbol W} x_i)
angle\ :\ \|a\|\leq 1ig\}$$

- Can't use random matrix theory
- $\sigma(Wx_i)$  has correlated entries
- Non-gaussian

#### Nonlinear features

$$egin{aligned} \hat{f}_a(x_i) &= ext{sign}(\langle a, u 
angle)\,, \ u_{ij} &= \sigma(\langle w_j, x_i 
angle) &= b_1 \langle w_j, x_i 
angle + b_* \sigma_\perp(\langle w_j, x_i 
angle) \end{aligned}$$

Noisy linear features

$$egin{aligned} \hat{f}_a(x_i) &= ext{sign}(\langle a, ilde{u} 
angle)\,, \ & ilde{u}_{ij} &= b_1 \langle w_j, x_i 
angle + b_* z_{ij}\,, \quad (z_{ij}) \sim_{iid} \mathsf{N}(0, 1) \end{aligned}$$

#### Nonlinear features

$$egin{aligned} \hat{f}_{m{a}}(m{x}_i) &= ext{sign}(\langlem{a},m{u}
angle)\,, \ &u_{ij} &= \sigma(\langlem{w}_j,m{x}_i
angle) &= b_1\langlem{w}_j,m{x}_i
angle + b_*\sigma_\perp(\langlem{w}_j,m{x}_i
angle) \end{aligned}$$

Noisy linear features

$$egin{aligned} \hat{f}_a(x_i) &= ext{sign}(\langle a, ilde{u} 
angle) \,, \ & ilde{u}_{ij} &= b_1 \langle w_j, x_i 
angle + b_* z_{ij} \,, \quad (z_{ij}) \sim_{iid} \mathsf{N}(0, 1) \end{aligned}$$

#### Nonlinear features

$$egin{aligned} \hat{f}_{m{a}}(m{x}_i) &= ext{sign}(\langlem{a},m{u}
angle)\,, \ u_{ij} &= \sigma(\langlem{w}_j,m{x}_i
angle) &= b_1\langlem{w}_j,m{x}_i
angle + b_*\sigma_\perp(\langlem{w}_j,m{x}_i
angle) \end{aligned}$$

Noisy linear features

$$egin{aligned} \hat{f}_a(x_i) &= ext{sign}(\langle a, ilde{u} 
angle) \,, \ & ilde{u}_{ij} &= b_1 \langle w_j, x_i 
angle + b_* z_{ij} \,, \quad (z_{ij}) \sim_{iid} \mathsf{N}(0, 1) \end{aligned}$$

#### Nonlinear features

$$egin{aligned} \hat{f}_{m{a}}(m{x}_i) &= ext{sign}(\langlem{a},m{u}
angle)\,, \ & u_{ij} &= \sigma(\langlem{w}_j,m{x}_i
angle) &= b_1\langlem{w}_j,m{x}_i
angle + b_*\sigma_\perp(\langlem{w}_j,m{x}_i
angle) \end{aligned}$$

#### Noisy linear features

$$egin{aligned} \hat{f}_{a}(x_{i}) &= ext{sign}(\langle a, ilde{u} 
angle)\,, \ & ilde{u}_{ij} &= b_{1}\langle w_{j}, x_{i} 
angle + b_{*}z_{ij}\,, \quad (z_{ij})\sim_{iid} \mathsf{N}(0,1) \end{aligned}$$

#### Nonlinear features

$$egin{aligned} \hat{f}_{m{a}}(m{x}_i) &= ext{sign}(\langlem{a},m{u}
angle)\,, \ u_{ij} &= \sigma(\langlem{w}_j,m{x}_i
angle) &= b_1\langlem{w}_j,m{x}_i
angle + b_*\sigma_\perp(\langlem{w}_j,m{x}_i
angle) \end{aligned}$$

#### Noisy linear features

$$egin{aligned} \hat{f}_{a}(x_{i}) &= ext{sign}(\langle a, ilde{u} 
angle)\,, \ & ilde{u}_{ij} &= b_{1}\langle w_{j}, x_{i} 
angle + b_{*}z_{ij}\,, \quad (z_{ij})\sim_{iid} \mathsf{N}(0,1) \end{aligned}$$

# Asymptotic equivalence

#### Theorem (Mei, M, 2019)

Consider ridge regression in the proportional asymptotics  $d \to \infty$ ,  $N/d \to \psi_1$ ,  $n/d \to \psi_2$ .

Then the nonlinear features model and the noisy linear features model are 'asymptotically equivalent:' they have asymptotically the same test error, training error, ...

#### Conjecture

Consider max-margin classification in the proportional asymptotics.

Then the nonlinear features model and the noisy linear features model are 'asymptotically equivalent.'

# Asymptotic equivalence

#### Theorem (Mei, M, 2019)

Consider ridge regression in the proportional asymptotics  $d \to \infty$ ,  $N/d \to \psi_1$ ,  $n/d \to \psi_2$ .

Then the nonlinear features model and the noisy linear features model are 'asymptotically equivalent:' they have asymptotically the same test error, training error, ...

#### Conjecture

Consider max-margin classification in the proportional asymptotics.

Then the nonlinear features model and the noisy linear features model are 'asymptotically equivalent.'

# Asymptotics of max-margin classification

Theorem (M, Ruan, Sogn, Yan, 2019)

Assume  $\mathbb{P}(y_i = +1 | x_i) = f_0(\langle \beta_0, x_i \rangle)$ ,  $x \sim \mathsf{N}(0, I_d)$ , and consider max-margin classification from noisy features  $\tilde{u}_{ij} = b_1 \langle w_j, x_i \rangle + b_* z_{ij}$ .

Let  $\kappa_n$  be the margin and  $\operatorname{Err}_n$  the test error. Then, we have

 $\mathrm{Err}_n = \mathrm{Err}_*(\zeta, \psi_1, \psi_2, f_0) + o_P(1), \quad \kappa_n = \kappa_*(\zeta, \psi_1, \psi_2, f_0) + o_P(1).$ 

where  $\mathrm{Err}_*(\zeta,\psi_1,\psi_2,f_0),\;\kappa_*(\zeta,\psi_1,\psi_2,f_0)$  are explicitly given.

▶ Paper: general Gaussian features.

▶ Statistical physics: Gardner 1988, Seung et al. 1992, ...

# Asymptotics of max-margin classification

Theorem (M, Ruan, Sogn, Yan, 2019)

Assume  $\mathbb{P}(y_i = +1 | x_i) = f_0(\langle \beta_0, x_i \rangle)$ ,  $x \sim \mathsf{N}(0, I_d)$ , and consider max-margin classification from noisy features  $\tilde{u}_{ij} = b_1 \langle w_j, x_i \rangle + b_* z_{ij}$ .

Let  $\kappa_n$  be the margin and  $\operatorname{Err}_n$  the test error. Then, we have

 $\mathrm{Err}_n = \mathrm{Err}_*(\zeta, \psi_1, \psi_2, f_0) + o_P(1), \quad \kappa_n = \kappa_*(\zeta, \psi_1, \psi_2, f_0) + o_P(1).$ 

where  $\operatorname{Err}_*(\zeta, \psi_1, \psi_2, f_0)$ ,  $\kappa_*(\zeta, \psi_1, \psi_2, f_0)$  are explicitly given.

▶ Paper: general Gaussian features.

▶ Statistical physics: Gardner 1988, Seung et al. 1992, ...

# Asymptotics of max-margin classification

Theorem (M, Ruan, Sogn, Yan, 2019)

Assume  $\mathbb{P}(y_i = +1|x_i) = f_0(\langle \beta_0, x_i \rangle)$ ,  $x \sim N(0, I_d)$ , and consider max-margin classification from noisy features  $\tilde{u}_{ij} = b_1 \langle w_j, x_i \rangle + b_* z_{ij}$ .

Let  $\kappa_n$  be the margin and  $\operatorname{Err}_n$  the test error. Then, we have

 $\mathrm{Err}_n = \mathrm{Err}_*(\zeta, \psi_1, \psi_2, f_0) + o_P(1), \quad \kappa_n = \kappa_*(\zeta, \psi_1, \psi_2, f_0) + o_P(1).$ 

where  $\operatorname{Err}_*(\zeta, \psi_1, \psi_2, f_0)$ ,  $\kappa_*(\zeta, \psi_1, \psi_2, f_0)$  are explicitly given.

- Paper: general Gaussian features.
- Statistical physics: Gardner 1988, Seung et al. 1992, ...

### Comparison



•  $\psi_2 = n/d = 2, d = 200, 400, \text{ ReLU}$  activations
### What is not captured by linear models?

What is not captured by linear models?

- Approximation
- Low-dimensional structures
- ► A 'simple' example

[see 1904.12191, 1906.08899]

What is not captured by linear models?

- Approximation
- Low-dimensional structures
- ► A 'simple' example

[see 1904.12191, 1906.08899]

### Setting

► 
$$x_i \sim \mathsf{N}(0, I_d),$$

$$y_i = f_*(x_i) \equiv b_0 + \langle x_i, Bx_i 
angle \; \; \; ext{ with } \; \; B \succeq 0.$$

• Optimum test error  $n = \infty$ 

 $R_{\mathsf{M},N} = \min_{\widehat{f}\in\mathcal{F}_{\mathsf{M},N}(oldsymbol{W})} \mathbb{E}\{(f_*(x) - \widehat{f}(x))^2\}, \quad \mathsf{M}\in\{\mathsf{RF},\mathsf{NT},\mathsf{NN}\}.$ 

• Here  $\sigma(x) = x^2$ 

(cf. paper for generalizations)

### Setting

► 
$$x_i \sim \mathsf{N}(0, I_d)$$
,

• Optimum test error  $n = \infty$ 

$$R_{\mathsf{M},N} = \min_{\widehat{f}\in\mathcal{F}_{\mathsf{M},N}(oldsymbol{W})} \mathbb{E}\{(f_*(oldsymbol{x}) - \widehat{f}(oldsymbol{x}))^2\}, \quad \mathsf{M}\in\{\mathsf{RF},\mathsf{NT},\mathsf{NN}\}.$$

• Here  $\sigma(x) = x^2$ 

(cf. paper for generalizations)

### Setting

• 
$$x_i \sim \mathsf{N}(0, I_d),$$
  
 $y_i = f_*(x_i) \equiv b_0 + \langle x_i, Bx_i 
angle ext{ with } B \succeq 0.$ 

• Optimum test error  $n = \infty$ 

$$R_{\mathsf{M},N} = \min_{\widehat{f}\in\mathcal{F}_{\mathsf{M},N}(oldsymbol{W})} \mathbb{E}\{(f_*(oldsymbol{x}) - \widehat{f}(oldsymbol{x}))^2\}, \quad \mathsf{M}\in\{\mathsf{RF},\mathsf{NT},\mathsf{NN}\}.$$

• Here  $\sigma(x) = x^2$ 

(cf. paper for generalizations)

### Experiment

- ▶  $m{B} \in \mathbb{R}^{450 imes 450}$ ,  $\lambda_i(m{b}) \sim_{iid} \exp(1)$
- N varies in  $\{30, \dots, 4500\}$ .



#### RF Model

 $f_{\sf RF}(x) = \sum_{i=1}^N a_i \sigma(\langle w_i, x \rangle)$  where  $w_i \sim {\sf N}(0,\Gamma)$  and  $a_i$  are trained.

Theorem (Ghorbani, Mei, Misiakiewicz, M., 2019)

Assume  $\sigma(x) = x^2 - 1$ . Then we have, as  $N, d \to \infty$  with  $N/d \to \psi_1$ :

$$R_{{\sf RF},N} = \|f_*\|_{L_2}^2 \left(1 - rac{\psi_1 d \langle m{B}, m{\Gamma} 
angle^2}{\|m{B}\|_F^2 (1 + \psi_1 d \|m{\Gamma}\|_F^2)} + o_{d,\mathbb{P}}(1)
ight)$$

#### • The correlation between $\Gamma$ and B controls the risk.

#### RF Model

 $f_{\mathsf{RF}}(x) = \sum_{i=1}^{N} a_i \sigma(\langle w_i, x \rangle)$  where  $w_i \sim \mathsf{N}(0, \Gamma)$  and  $a_i$  are trained.

Theorem (Ghorbani, Mei, Misiakiewicz, M., 2019)

Assume  $\sigma(x) = x^2 - 1$ . Then we have, as  $N, d \to \infty$  with  $N/d \to \psi_1$ :

$$R_{\mathsf{RF},N} = \|f_*\|_{L_2}^2 \left(1 - rac{\psi_1 d \langle m{B}, m{\Gamma} 
angle^2}{\|m{B}\|_F^2 (1 + \psi_1 d \|m{\Gamma}\|_F^2)} + o_{d,\mathbb{P}}(1)
ight)\,.$$

• The correlation between  $\Gamma$  and B controls the risk.

#### RF Model

 $f_{\sf RF}(x) = \sum_{i=1}^N a_i \sigma(\langle w_i, x \rangle)$  where  $w_i \sim {\sf N}(0,\Gamma)$  and  $a_i$  are trained.

Theorem (Ghorbani, Mei, Misiakiewicz, M., 2019)

Assume  $\sigma(x) = x^2 - 1$ . Then we have, as  $N, d \to \infty$  with  $N/d \to \psi_1$ :

$$R_{\mathsf{RF},N} = \|f_*\|_{L_2}^2 \left(1 - rac{\psi_1 d \langle m{B}, m{\Gamma} 
angle^2}{\|m{B}\|_F^2 (1 + \psi_1 d \|m{\Gamma}\|_F^2)} + o_{d,\mathbb{P}}(1)
ight)\,.$$

• The correlation between  $\Gamma$  and B controls the risk.

$$\lim_{\psi_1 o \infty} \lim_{d o \infty, N/d o \psi_1} rac{R_{\mathsf{F},N}}{\mathbb{E}\{f^2_*\}} = \lim_{d o \infty} \left(1 - rac{\langle m{\Gamma}, m{B} 
angle^2}{\|m{\Gamma}\|_F^2\|m{B}\|_F^2}
ight) \,.$$

- Risk vanishes only if  $\Gamma$  is chosen perfectly and  $\psi_1 \to \infty$ .
- ▶ This result is true for any activation
- ▶ The asymptotic risk is independent of the non-linearity!

$$\lim_{\psi_1 o\infty} \lim_{d o\infty,N/d o\psi_1} rac{R_{\mathsf{RF},N}}{\mathbb{E}\{f^2_*\}} = \lim_{d o\infty} \left(1 - rac{\langle m{\Gamma},m{B}
angle^2}{\|m{\Gamma}\|_F^2\|m{B}\|_F^2}
ight)\,.$$

- Risk vanishes only if  $\Gamma$  is chosen perfectly and  $\psi_1 o \infty$ .
- ▶ This result is true for any activation
- ▶ The asymptotic risk is independent of the non-linearity!

$$\lim_{\psi_1 o \infty} \lim_{d o \infty, N/d o \psi_1} rac{R_{\mathsf{RF},N}}{\mathbb{E}\{f^2_*\}} = \lim_{d o \infty} \left(1 - rac{\langle m{\Gamma}, m{B} 
angle^2}{\|m{\Gamma}\|_F^2 \|m{B}\|_F^2}
ight).$$

- Risk vanishes only if  $\Gamma$  is chosen perfectly and  $\psi_1 o \infty$ .
- ▶ This result is true for any activation
- ▶ The asymptotic risk is independent of the non-linearity!

٠

$$\lim_{\psi_1 o \infty} \lim_{d o \infty, N/d o \psi_1} rac{R_{\mathsf{RF},N}}{\mathbb{E}\{f^2_*\}} = \lim_{d o \infty} \left(1 - rac{\langle m{\Gamma}, m{B} 
angle^2}{\|m{\Gamma}\|_F^2 \|m{B}\|_F^2}
ight).$$

- Risk vanishes only if  $\Gamma$  is chosen perfectly and  $\psi_1 o \infty$ .
- ▶ This result is true for any activation
- ▶ The asymptotic risk is independent of the non-linearity!

.

▶ Naive RF does not learn an efficient representation of the data.



▶ Naive RF does not learn an efficient representation of the data.



### NT Model

•  $f_{\mathsf{NT}}(x) = c + \sum_{i=1}^N \sigma'(\langle w_i, x \rangle) \langle a_i, x \rangle$  where  $w_i \sim_{i.i.d} \mathsf{N}(0, I_d/d)$ .

Theorem (Ghorbani, Mei, Misiakiewicz, M., 2019)

For  $N, d 
ightarrow \infty$  with  $N/d 
ightarrow \psi_1$ 

$$\frac{\mathbb{E}[R_{\mathsf{NT},N}]}{\mathbb{E}\{f_*^2\}} = \Big\{(1-\psi_1)_+^2 + \psi_1(1-\psi_1)_+ \frac{\mathrm{Tr}(B)^2}{d\|B\|_F^2} + o_d(1)\Big\}$$

where the expectation is taken over  $w_i \sim_{i,i,d} N(\mathbf{0}, \mathbf{I}_d/d)$ .

### NT Model

•  $f_{\mathsf{NT}}(x) = c + \sum_{i=1}^N \sigma'(\langle w_i, x \rangle) \langle a_i, x \rangle$  where  $w_i \sim_{i.i.d} \mathsf{N}(0, I_d/d)$ .

Theorem (Ghorbani, Mei, Misiakiewicz, M., 2019)

For N,  $d 
ightarrow \infty$  with  $N/d 
ightarrow \psi_1$ 

$$rac{\mathbb{E}[R_{\mathsf{NT},N}]}{\mathbb{E}\{f_*^2\}} = \Big\{(1-\psi_1)_+^2 + \psi_1(1-\psi_1)_+ rac{\mathrm{Tr}(B)^2}{d\|B\|_F^2} + o_d(1)\Big\}$$

where the expectation is taken over  $w_i \sim_{i.i.d} N(0, I_d/d)$ .

Quadratic Model - NT



Theorem (Ghorbani, Mei, Misiakiewicz, M., 2019) For  $N, d \to \infty$  with  $N/d \to \psi_1$  $\frac{\mathbb{E}[R_{\mathsf{NT},N]}}{\mathbb{E}\{f_*^2\}} = \left\{ (1 - \psi_1)_+^2 + \psi_1(1 - \psi_1)_+ \frac{\operatorname{Tr}(B)^2}{d||B||_F^2} + o_d(1) \right\}$ where the expectation is taken over  $w_i \sim_{i,i,d} \mathsf{N}(\mathbf{0}, \mathbf{I}_d/d)$ .

• Does NT learn subspaces based on their importance in B?

**No!** NT fits random directions (but more parameters).

Theorem (Ghorbani, Mei, Misiakiewicz, M., 2019) For  $N, d \to \infty$  with  $N/d \to \psi_1$  $\frac{\mathbb{E}[R_{\mathsf{NT},N]}}{\mathbb{E}\{f_*^2\}} = \left\{ (1 - \psi_1)_+^2 + \psi_1(1 - \psi_1)_+ \frac{\operatorname{Tr}(B)^2}{d||B||_F^2} + o_d(1) \right\}$ where the expectation is taken over  $w_i \sim_{i,i,d} \mathsf{N}(\mathbf{0}, \mathbf{I}_d/d)$ .

- ▶ Does NT learn subspaces based on their importance in B?
- **No!** NT fits random directions (but more parameters).

Theorem (Ghorbani, Mei, Misiakiewicz, M., 2019) For  $N, d \to \infty$  with  $N/d \to \psi_1$  $\frac{\mathbb{E}[R_{\mathsf{NT},N]}}{\mathbb{E}\{f_*^2\}} = \left\{ (1 - \psi_1)_+^2 + \psi_1(1 - \psi_1)_+ \frac{\operatorname{Tr}(B)^2}{d \|B\|_F^2} + o_d(1) \right\}$ where the expectation is taken over  $w_i \sim_{i,i,d} \mathsf{N}(\mathbf{0}, \mathbf{I}_d/d)$ .

- ▶ Does NT learn subspaces based on their importance in B?
- ▶ No! NT fits random directions (but more parameters).

- Square non-linearity  $\longrightarrow f_{NN}(x; W, c) = \sum_{i=1}^{N} \langle w_i, x \rangle^2 + c.$
- $\blacktriangleright \ R(f_{\mathsf{NN}}) = R(\boldsymbol{W}, \boldsymbol{c}) = \mathbb{E}\Big[\Big(\langle \boldsymbol{x}\boldsymbol{x}^{\mathsf{T}}, \boldsymbol{B} \boldsymbol{W}\boldsymbol{W}^{\mathsf{T}}\rangle + b_0 \boldsymbol{c}\Big)^2\Big].$

#### Does gradient descent converge to this value?

- Square non-linearity  $\longrightarrow f_{NN}(x; W, c) = \sum_{i=1}^{N} \langle w_i, x \rangle^2 + c.$
- $\blacktriangleright \ R(f_{\mathsf{NN}}) = R(\boldsymbol{W}, \boldsymbol{c}) = \mathbb{E}\Big[\Big(\langle \boldsymbol{x}\boldsymbol{x}^{\mathsf{T}}, \boldsymbol{B} \boldsymbol{W}\boldsymbol{W}^{\mathsf{T}}\rangle + b_0 \boldsymbol{c}\Big)^2\Big].$

#### Does gradient descent converge to this value?

• Square non-linearity 
$$\longrightarrow f_{\sf NN}(x; W, c) = \sum_{i=1}^N \langle w_i, x \rangle^2 + c.$$

$$\blacktriangleright \ R(f_{\mathsf{NN}}) = R(W, c) = \mathbb{E}\Big[\Big(\langle xx^{\mathsf{T}}, B - WW^{\mathsf{T}} \rangle + b_0 - c\Big)^2\Big].$$

### ▶ Does gradient descent converge to this value?

• Square non-linearity  $\longrightarrow f_{\sf NN}(x; W, c) = \sum_{i=1}^N \langle w_i, x \rangle^2 + c.$ 

$$\blacktriangleright \ R(f_{\mathsf{NN}}) = R(\boldsymbol{W}, \boldsymbol{c}) = \mathbb{E}\Big[\Big(\langle \boldsymbol{x}\boldsymbol{x}^{\mathsf{T}}, \boldsymbol{B} - \boldsymbol{W}\boldsymbol{W}^{\mathsf{T}}\rangle + b_0 - \boldsymbol{c}\Big)^2\Big].$$

#### ▶ Does gradient descent converge to this value?

• Square non-linearity  $\longrightarrow f_{\sf NN}(x; W, c) = \sum_{i=1}^N \langle w_i, x \rangle^2 + c.$ 

$$\blacktriangleright \ R(f_{\mathsf{NN}}) = R(\boldsymbol{W}, \boldsymbol{c}) = \mathbb{E}\Big[\Big(\langle \boldsymbol{x}\boldsymbol{x}^{\mathsf{T}}, \boldsymbol{B} - \boldsymbol{W}\boldsymbol{W}^{\mathsf{T}}\rangle + b_0 - \boldsymbol{c}\Big)^2\Big].$$

#### Does gradient descent converge to this value?

### NN with Gradient Descent

One-pass version of SGD,

$$egin{aligned} & (m{W}_{k+1}, m{c}_{k+1}) = (m{W}_k, m{c}_k) - m{arepsilon} 
abla_{m{W}, m{c}} \Big( f_*(m{x}_k) - \hat{f}(m{x}_k; m{W}, m{c}) \Big)^2 \, . \end{aligned}$$

▶ Define

$$R(\ell,\varepsilon)\equiv R(\boldsymbol{W}_{\ell},c_{\ell}).$$

Theorem (Ghorbani, Mei, Misiakiewicz, M., 2019)

We have

$$\lim_{t o\infty}\lim_{arepsilon o \infty} \mathbb{P}ig( \Big| R(\ell=t/arepsilon,arepsilon) - \inf_{oldsymbol{W},c} R(oldsymbol{W},c) \Big| \geq \delta ) = 0,$$

(probability is over the initialization  $(W_0, c_0)$  and the samples.)

### NN with Gradient Descent

One-pass version of SGD,

$$( \, oldsymbol{W}_{k+1}, \, oldsymbol{c}_{k+1}) = ( \, oldsymbol{W}_k, \, oldsymbol{c}_k) - arepsilon 
abla _{oldsymbol{W}, \, c} \Big( f_*(oldsymbol{x}_k) - \hat{f}(oldsymbol{x}_k; \, oldsymbol{W}, \, c) \Big)^2$$

Define

$$R(\ell, \varepsilon) \equiv R(W_{\ell}, c_{\ell}).$$

#### Theorem (Ghorbani, Mei, Misiakiewicz, M., 2019)

We have

$$\lim_{t o\infty}\lim_{arepsilon o \infty} \mathbb{P}ig( \Big| R(\ell=t/arepsilon,arepsilon) - \inf_{oldsymbol{W},c} R(oldsymbol{W},c) \Big| \geq \delta ) = 0,$$

(probability is over the initialization  $(W_0, c_0)$  and the samples.)

### NN with Gradient Descent

One-pass version of SGD,

$$( \, oldsymbol{W}_{k+1}, \, oldsymbol{c}_{k+1}) = ( \, oldsymbol{W}_k, \, oldsymbol{c}_k) - arepsilon 
abla _{oldsymbol{W}, \, c} \Big( f_*(oldsymbol{x}_k) - \hat{f}(oldsymbol{x}_k; \, oldsymbol{W}, \, c) \Big)^2$$

Define

$$R(\ell, \varepsilon) \equiv R(W_{\ell}, c_{\ell}).$$

Theorem (Ghorbani, Mei, Misiakiewicz, M., 2019)

We have

$$\lim_{\epsilon o \infty} \lim_{arepsilon o 0} \mathbb{P}ig( \Big| R(oldsymbol{\ell} = t/arepsilon, arepsilon) - \inf_{oldsymbol{W}, c} R(oldsymbol{W}, c) \Big| \geq \delta ig) = 0,$$

(probability is over the initialization  $(W_0, c_0)$  and the samples.)



$$rac{R_{\mathsf{M},N}}{\mathbb{E}\{f_*^2\}} pprox \left\{ egin{array}{ll} 1 - rac{\psi_1}{1 + \psi_1} rac{\mathrm{Tr}(m{B})^2}{d\|m{B}\|_F^2} & ext{for } \mathsf{M} = \mathsf{RF}, \ (1 - \psi_1)_+^2 + \psi_1(1 - \psi_1)_+ rac{\mathrm{Tr}(m{B})^2}{d\|m{B}\|_F^2} & ext{for } \mathsf{M} = \mathsf{NT}, \ 1 - rac{\sum_{i=1}^{d \wedge N} \lambda_i(m{B})^2}{\|m{B}\|_F^2} & ext{for } \mathsf{M} = \mathsf{NN}. \end{array} 
ight.$$

- $\blacktriangleright \exists B$  arbtrarily large gap between NN and NT.
- ▶ How general are these phenomena?

$$rac{R_{\mathsf{M},N}}{\mathbb{E}\{f_*^2\}} pprox \left\{ egin{array}{ll} 1 - rac{\psi_1}{1 + \psi_1} rac{\mathrm{Tr}(B)^2}{d\|B\|_F^2} & ext{for } \mathsf{M} = \mathsf{RF}, \ (1 - \psi_1)_+^2 + \psi_1(1 - \psi_1)_+ rac{\mathrm{Tr}(B)^2}{d\|B\|_F^2} & ext{for } \mathsf{M} = \mathsf{NT}, \ 1 - rac{\sum_{i=1}^{d \wedge N} \lambda_i(B)^2}{\|B\|_F^2} & ext{for } \mathsf{M} = \mathsf{NN}. \end{array} 
ight.$$

- $\blacktriangleright \exists B$  arbtrarily large gap between NN and NT.
- ▶ How general are these phenomena?

$$rac{R_{\mathsf{M},N}}{\mathbb{E}\{f_*^2\}} pprox \left\{ egin{array}{ll} 1 - rac{\psi_1}{1 + \psi_1} rac{\mathrm{Tr}(B)^2}{d\|B\|_F^2} & ext{for } \mathsf{M} = \mathsf{RF}, \ (1 - \psi_1)_+^2 + \psi_1(1 - \psi_1)_+ rac{\mathrm{Tr}(B)^2}{d\|B\|_F^2} & ext{for } \mathsf{M} = \mathsf{NT}, \ 1 - rac{\sum_{i=1}^{d \wedge N} \lambda_i(B)^2}{\|B\|_F^2} & ext{for } \mathsf{M} = \mathsf{NN}. \end{array} 
ight.$$

- $\blacktriangleright \exists B$  arbtrarily large gap between NN and NT.
- ▶ How general are these phenomena?

$$rac{R_{\mathsf{M},N}}{\mathbb{E}\{f_*^2\}} pprox \left\{ egin{array}{ll} 1 - rac{\psi_1}{1 + \psi_1} rac{\mathrm{Tr}(B)^2}{d\|B\|_F^2} & ext{for } \mathsf{M} = \mathsf{RF}, \ (1 - \psi_1)_+^2 + \psi_1(1 - \psi_1)_+ rac{\mathrm{Tr}(B)^2}{d\|B\|_F^2} & ext{for } \mathsf{M} = \mathsf{NT}, \ 1 - rac{\sum_{i=1}^{d \wedge N} \lambda_i(B)^2}{\|B\|_F^2} & ext{for } \mathsf{M} = \mathsf{NN}. \end{array} 
ight.$$

- $\blacktriangleright \exists B$  arbtrarily large gap between NN and NT.
- ▶ How general are these phenomena?

$$rac{R_{\mathsf{M},N}}{\mathbb{E}\{f_*^2\}} pprox \left\{ egin{array}{ll} 1 - rac{\psi_1}{1+\psi_1}rac{\mathrm{Tr}(m{B})^2}{d\|m{B}\|_F^2} & ext{for } \mathsf{M} = \mathsf{RF}, \ (1-\psi_1)_+^2 + \psi_1(1-\psi_1)_+rac{\mathrm{Tr}(m{B})^2}{d\|m{B}\|_F^2} & ext{for } \mathsf{M} = \mathsf{NT}, \ 1 - rac{\sum_{i=1}^{d\wedge N}\lambda_i(m{B})^2}{\|m{B}\|_F^2} & ext{for } \mathsf{M} = \mathsf{NN}. \end{array} 
ight.$$

- ▶ Naive RF, NT do not learn good representations of the data.
- ▶  $\exists B$  arbtrarily large gap between NN and NT.
- ▶ How general are these phenomena?
## Comparison

$$rac{R_{\mathsf{M},N}}{\mathbb{E}\{f_*^2\}} pprox \left\{ egin{array}{ll} 1 - rac{\psi_1}{1+\psi_1}rac{\mathrm{Tr}(m{B})^2}{d\|m{B}\|_F^2} & ext{for } \mathsf{M} = \mathsf{RF}, \ (1-\psi_1)_+^2 + \psi_1(1-\psi_1)_+rac{\mathrm{Tr}(m{B})^2}{d\|m{B}\|_F^2} & ext{for } \mathsf{M} = \mathsf{NT}, \ 1 - rac{\sum_{i=1}^{d\wedge N}\lambda_i(m{B})^2}{\|m{B}\|_F^2} & ext{for } \mathsf{M} = \mathsf{NN}. \end{array} 
ight.$$

- ▶ Naive RF, NT do not learn good representations of the data.
- ▶  $\exists B$  arbtrarily large gap between NN and NT.
- How general are these phenomena?

# Comparison with ReLU



#### Conclusion

## Conclusion

#### ▶ Thermodynamic/wide limit: a useful concept

Different regimes (depends on training)

▶ Linearized regime captures certain phenomena

Thanks!

### Conclusion

- ▶ Thermodynamic/wide limit: a useful concept
- Different regimes (depends on training)
- Linearized regime captures certain phenomena

Thanks!