

# Limitations of Lazy Training of Two-layers Neural Networks

Theodor Misiakiewicz

Stanford University

December 11, 2019

Joint work with Behrooz Ghorbani, Song Mei, Andrea Montanari

# Models

- ▶ Two-layers Neural Network (NN) model:

$$f_{\text{NN}}(\mathbf{x}; \mathbf{a}, \mathbf{W}) = \sum_{i=1}^N a_i \sigma(\mathbf{w}_i^\top \mathbf{x})$$

*Lazy training regime or Kernel limit* [Du et al.,18], [Chizat, Bach,18]

- ▶ Random Features (RF) model: (second layer linearization)

$$f_{\text{RF}}(\mathbf{x}; \mathbf{a}) = \sum_{i=1}^N a_i \sigma(\mathbf{w}_i^\top \mathbf{x}) \quad [\text{Rahimi, Recht, 08}]$$

- ▶ Neural Tangent (NT) model: (first layer linearization)

$$f_{\text{NT}}(\mathbf{x}; \mathbf{a}) = \sum_{i=1}^N \langle \mathbf{a}_i, \mathbf{x} \rangle \sigma'(\mathbf{w}_i^\top \mathbf{x}) \quad [\text{Jacot et al., 18}]$$

# Models

- ▶ **Two-layers Neural Network (NN) model:**

$$f_{\text{NN}}(\mathbf{x}; \mathbf{a}, \mathbf{W}) = \sum_{i=1}^N a_i \sigma(\mathbf{w}_i^\top \mathbf{x})$$

*Lazy training regime or Kernel limit [Du et al.,18], [Chizat, Bach,18]*

- ▶ **Random Features (RF) model:** (second layer linearization)

$$f_{\text{RF}}(\mathbf{x}; \mathbf{a}) = \sum_{i=1}^N a_i \sigma(\mathbf{w}_i^\top \mathbf{x}) \quad [\text{Rahimi, Recht, 08}]$$

- ▶ **Neural Tangent (NT) model:** (first layer linearization)

$$f_{\text{NT}}(\mathbf{x}; \mathbf{a}) = \sum_{i=1}^N \langle \mathbf{a}_i, \mathbf{x} \rangle \sigma'(\mathbf{w}_i^\top \mathbf{x}) \quad [\text{Jacot et al., 18}]$$

# Models

- ▶ **Two-layers Neural Network (NN) model:**

$$f_{\text{NN}}(\mathbf{x}; \mathbf{a}, \mathbf{W}) = \sum_{i=1}^N a_i \sigma(\mathbf{w}_i^\top \mathbf{x})$$

*Lazy training regime* or *Kernel limit* [Du et al.,18], [Chizat, Bach,18]

- ▶ **Random Features (RF) model:** (second layer linearization)

$$f_{\text{RF}}(\mathbf{x}; \mathbf{a}) = \sum_{i=1}^N a_i \sigma(\mathbf{w}_i^\top \mathbf{x}) \quad [\text{Rahimi, Recht, 08}]$$

- ▶ **Neural Tangent (NT) model:** (first layer linearization)

$$f_{\text{NT}}(\mathbf{x}; \mathbf{a}) = \sum_{i=1}^N \langle \mathbf{a}_i, \mathbf{x} \rangle \sigma'(\mathbf{w}_i^\top \mathbf{x}) \quad [\text{Jacot et al., 18}]$$

# Models

- ▶ **Two-layers Neural Network (NN) model:**

$$f_{\text{NN}}(\mathbf{x}; \mathbf{a}, \mathbf{W}) = \sum_{i=1}^N a_i \sigma(\mathbf{w}_i^\top \mathbf{x})$$

*Lazy training regime* or *Kernel limit* [Du et al.,18], [Chizat, Bach,18]

- ▶ **Random Features (RF) model:** (second layer linearization)

$$f_{\text{RF}}(\mathbf{x}; \mathbf{a}) = \sum_{i=1}^N a_i \sigma(\mathbf{w}_i^\top \mathbf{x}) \quad [\text{Rahimi, Recht, 08}]$$

- ▶ **Neural Tangent (NT) model:** (first layer linearization)

$$f_{\text{NT}}(\mathbf{x}; \mathbf{a}) = \sum_{i=1}^N \langle \mathbf{a}_i, \mathbf{x} \rangle \sigma'(\mathbf{w}_i^\top \mathbf{x}) \quad [\text{Jacot et al., 18}]$$

# Models

- ▶ **Two-layers Neural Network (NN) model:**

$$f_{\text{NN}}(\mathbf{x}; \mathbf{a}, \mathbf{W}) = \sum_{i=1}^N \mathbf{a}_i \sigma(\mathbf{w}_i^\top \mathbf{x})$$

*Lazy training regime* or *Kernel limit* [Du et al.,18], [Chizat, Bach,18]

- ▶ **Random Features (RF) model:** (second layer linearization)

$$f_{\text{RF}}(\mathbf{x}; \mathbf{a}) = \sum_{i=1}^N \mathbf{a}_i \sigma(\mathbf{w}_i^\top \mathbf{x}) \quad [\text{Rahimi, Recht, 08}]$$

- ▶ **Neural Tangent (NT) model:** (first layer linearization)

$$f_{\text{NT}}(\mathbf{x}; \mathbf{a}) = \sum_{i=1}^N \langle \mathbf{a}_i, \mathbf{x} \rangle \sigma'(\mathbf{w}_i^\top \mathbf{x}) \quad [\text{Jacot et al., 18}]$$

## Goal: compare these three models.

### Setting:

▶  $\mathbf{x}_i \sim \mathcal{N}(0, \mathbf{I}_d)$ ,

$$y_i = f_*(\mathbf{x}_i) \equiv \langle \mathbf{x}_i, \mathbf{B}\mathbf{x}_i \rangle + b_0, \quad \text{with } B \succeq 0$$

▶ Activation function  $\sigma(x) = x^2$

▶  $N$  number of neurons:  $N/d = \rho$  and  $d$  large (high-dimensional regime)

We compare the population squared loss:

$$R_{M,N}(f^*) = \min_{\hat{f} \in \mathcal{F}_M} \mathbb{E}_{\mathbf{x}} \left\{ \left( f_*(\mathbf{x}) - \hat{f}(\mathbf{x}) \right)^2 \right\}, \quad M \in \{\text{RF}, \text{NT}\}$$

$$R_{\text{NN},N}(f^*) = \lim_{t \rightarrow \infty} \lim_{\varepsilon \rightarrow 0} \mathbb{E}_{\mathbf{x}} \left\{ \left( f_*(\mathbf{x}) - \hat{f}_{\text{SGD}}(\mathbf{x}; t/\varepsilon, \varepsilon) \right)^2 \right\}$$

## Goal: compare these three models.

### Setting:

▶  $\mathbf{x}_i \sim \mathcal{N}(0, \mathbf{I}_d)$ ,

$$y_i = f_*(\mathbf{x}_i) \equiv \langle \mathbf{x}_i, \mathbf{B}\mathbf{x}_i \rangle + b_0, \quad \text{with } \mathbf{B} \succeq 0$$

▶ Activation function  $\sigma(x) = x^2$

▶  $N$  number of neurons:  $N/d = \rho$  and  $d$  large (high-dimensional regime)

We compare the population squared loss:

$$R_{M,N}(f^*) = \min_{\hat{f} \in \mathcal{F}_M} \mathbb{E}_{\mathbf{x}} \left\{ \left( f_*(\mathbf{x}) - \hat{f}(\mathbf{x}) \right)^2 \right\}, \quad M \in \{\text{RF, NT}\}$$

$$R_{\text{NN},N}(f^*) = \lim_{t \rightarrow \infty} \lim_{\varepsilon \rightarrow 0} \mathbb{E}_{\mathbf{x}} \left\{ \left( f_*(\mathbf{x}) - \hat{f}_{\text{SGD}}(\mathbf{x}; t/\varepsilon, \varepsilon) \right)^2 \right\}$$



## Goal: compare these three models.

### Setting:

▶  $\mathbf{x}_i \sim \mathcal{N}(0, \mathbf{I}_d)$ ,

$$y_i = f_*(\mathbf{x}_i) \equiv \langle \mathbf{x}_i, \mathbf{B}\mathbf{x}_i \rangle + b_0, \quad \text{with } \mathbf{B} \succeq 0$$

▶ Activation function  $\sigma(x) = x^2$

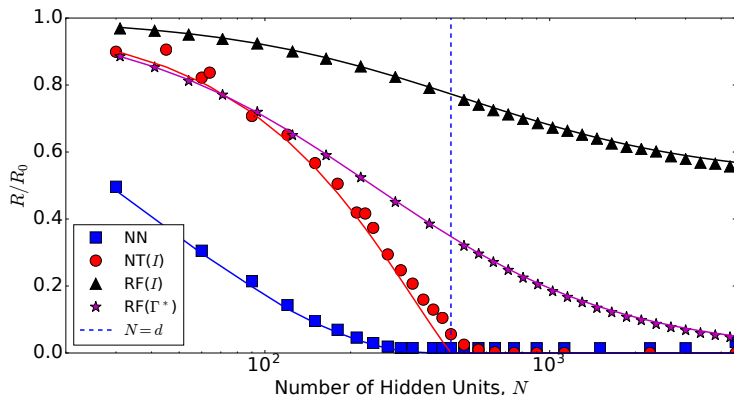
▶  $N$  number of neurons:  $N/d = \rho$  and  $d$  large (high-dimensional regime)

We compare the population squared loss:

$$R_{M,N}(f^*) = \min_{\hat{f} \in \mathcal{F}_M} \mathbb{E}_{\mathbf{x}} \left\{ \left( f_*(\mathbf{x}) - \hat{f}(\mathbf{x}) \right)^2 \right\}, \quad M \in \{\text{RF}, \text{NT}\}$$

$$R_{\text{NN},N}(f^*) = \lim_{t \rightarrow \infty} \lim_{\varepsilon \rightarrow 0} \mathbb{E}_{\mathbf{x}} \left\{ \left( f_*(\mathbf{x}) - \hat{f}_{\text{SGD}}(\mathbf{x}; t/\varepsilon, \varepsilon) \right)^2 \right\}$$

# Results



**Figure:** Population error in fitting a quadratic function in  $d = 450$  dimensions for random features (RF), neural tangent (NT), and SGD trained neural networks (NN). Lines are analytical predictions and dots are empirical results.

# RF vs NT vs NN

- ▶ RF model does not capture quadratic functions
- ▶ NN model:

$$R_{\text{NN},N}(f_*) = \min_{\mathbf{W} \in \mathbb{R}^{N \times d}} 2 \|\mathbf{B} - \mathbf{W}\mathbf{W}^\top\|_F^2 = 2 \sum_{i=N \wedge d + 1}^d \lambda_i(\mathbf{B})^2$$

with  $\lambda_1(\mathbf{B}) \geq \lambda_2(\mathbf{B}) \geq \dots \geq \lambda_d(\mathbf{B}) \geq 0$ .

- ▶ NT model:  $\mathbf{W} = [\mathbf{w}_1^\top, \dots, \mathbf{w}_N^\top]^\top \in \mathbb{R}^{N \times d}$  fixed at initialization

$$R_{\text{NT},N}(f_*) = \min_{\mathbf{A} \in \mathbb{R}^{N \times d}} 2 \|\mathbf{B} - \mathbf{W}\mathbf{A}^\top - \mathbf{A}\mathbf{W}^\top\|_F^2$$

for  $N < d$ , fit along  $\mathbf{W}$  random subspace of dimension  $N$ .

# RF vs NT vs NN

- ▶ RF model does not capture quadratic functions
- ▶ NN model:

$$R_{\text{NN},N}(f_*) = \min_{\mathbf{W} \in \mathbb{R}^{N \times d}} 2 \|\mathbf{B} - \mathbf{W}\mathbf{W}^\top\|_F^2 = 2 \sum_{i=N \wedge d + 1}^d \lambda_i(\mathbf{B})^2$$

with  $\lambda_1(\mathbf{B}) \geq \lambda_2(\mathbf{B}) \geq \dots \geq \lambda_d(\mathbf{B}) \geq 0$ .

- ▶ NT model:  $\mathbf{W} = [\mathbf{w}_1^\top, \dots, \mathbf{w}_N^\top]^\top \in \mathbb{R}^{N \times d}$  fixed at initialization

$$R_{\text{NT},N}(f_*) = \min_{\mathbf{A} \in \mathbb{R}^{N \times d}} 2 \|\mathbf{B} - \mathbf{W}\mathbf{A}^\top - \mathbf{A}\mathbf{W}^\top\|_F^2$$

for  $N < d$ , fit along  $\mathbf{W}$  random subspace of dimension  $N$ .

# RF vs NT vs NN

- ▶ RF model does not capture quadratic functions
- ▶ NN model:

$$R_{\text{NN},N}(f_*) = \min_{\mathbf{W} \in \mathbb{R}^{N \times d}} 2 \|\mathbf{B} - \mathbf{W}\mathbf{W}^\top\|_F^2 = 2 \sum_{i=N \wedge d + 1}^d \lambda_i(\mathbf{B})^2$$

with  $\lambda_1(\mathbf{B}) \geq \lambda_2(\mathbf{B}) \geq \dots \geq \lambda_d(\mathbf{B}) \geq 0$ .

- ▶ NT model:  $\mathbf{W} = [\mathbf{w}_1^\top, \dots, \mathbf{w}_N^\top]^\top \in \mathbb{R}^{N \times d}$  fixed at initialization

$$R_{\text{NT},N}(f_*) = \min_{\mathbf{A} \in \mathbb{R}^{N \times d}} 2 \|\mathbf{B} - \mathbf{W}\mathbf{A}^\top - \mathbf{A}\mathbf{W}^\top\|_F^2$$

for  $N < d$ , fit along  $\mathbf{W}$  random subspace of dimension  $N$ .

# RF vs NT vs NN

- ▶ RF model does not capture quadratic functions
- ▶ NN model:

$$R_{\text{NN},N}(f_*) = \min_{\mathbf{W} \in \mathbb{R}^{N \times d}} 2 \|\mathbf{B} - \mathbf{W}\mathbf{W}^\top\|_F^2 = 2 \sum_{i=N \wedge d + 1}^d \lambda_i(\mathbf{B})^2$$

with  $\lambda_1(\mathbf{B}) \geq \lambda_2(\mathbf{B}) \geq \dots \geq \lambda_d(\mathbf{B}) \geq 0$ .

- ▶ NT model:  $\mathbf{W} = [\mathbf{w}_1^\top, \dots, \mathbf{w}_N^\top]^\top \in \mathbb{R}^{N \times d}$  fixed at initialization

$$R_{\text{NT},N}(f_*) = \min_{\mathbf{A} \in \mathbb{R}^{N \times d}} 2 \|\mathbf{B} - \mathbf{W}\mathbf{A}^\top - \mathbf{A}\mathbf{W}^\top\|_F^2$$

for  $N < d$ , fit along  $\mathbf{W}$  random subspace of dimension  $N$ .

# Interpretation

- ▶ Fully trained NN learns the most important eigendirections, while the NT model remains confined to a random set of directions.

Neural networks are superior to linearized model such as RF and NT, because they can learn a good representation of the data

Mechanism more general: mixture of Gaussians, ReLu activation...

# Interpretation

- ▶ Fully trained NN learns the most important eigendirections, while the NT model remains confined to a random set of directions.

Neural networks are superior to linearized model such as RF and NT, because they can learn a good representation of the data

Mechanism more general: mixture of Gaussians, ReLu activation...



# Interpretation

- ▶ Fully trained NN learns the most important eigendirections, while the NT model remains confined to a random set of directions.

Neural networks are superior to linearized model such as RF and NT, because they can learn a good representation of the data

Mechanism more general: mixture of Gaussians, ReLu activation...

# Thank you!

For further discussions, you can visit our poster:

**Poster # 230**  
**East Exhibition Hall B + C**  
**5:00 - 7:00pm, Wednesday 11th**

If you have any questions: please email us at [misiakie@stanford.edu](mailto:misiakie@stanford.edu)