

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

*Equal contributions

Lazy Training Regime

Two-layers Neural Network (NN):

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

Lazy Training regime: for some initialization $\theta_i^0 = (a_i^0, \mathbf{w}_i^0)$

$$\begin{aligned} f_{\text{NN}}(\mathbf{x}; \theta^t) &\approx f_{\text{NN}}(\mathbf{x}; \theta^0) + \langle \theta^t - \theta^0, \nabla_{\theta} f_{\text{NN}}(\mathbf{x}; \theta^0) \rangle \\ &\approx 0 + \underbrace{\sum_{i=1}^N t_i \sigma(\langle \mathbf{w}_i^0, \mathbf{x} \rangle)}_{\text{Random Features (RF) model}} + \underbrace{\sum_{i=1}^N \langle b_i, \mathbf{x} \rangle \sigma'(\langle \mathbf{w}_i^0, \mathbf{x} \rangle)}_{\text{Neural Tangent (NT) model}} \end{aligned}$$

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

Lazy Training Regime

Two-layers Neural Network (NN):

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

Lazy Training regime: for some initialization $\theta_i^0 = (a_i^0, \mathbf{w}_i^0)$

$$\begin{aligned} f_{\text{NN}}(\mathbf{x}; \theta^t) &\approx f_{\text{NN}}(\mathbf{x}; \theta^0) + \langle \theta^t - \theta^0, \nabla_{\theta} f_{\text{NN}}(\mathbf{x}; \theta^0) \rangle \\ &\approx 0 + \underbrace{\sum_{i=1}^N t_i \sigma(\langle \mathbf{w}_i^0, \mathbf{x} \rangle)}_{\text{Random Features (RF) model}} + \underbrace{\sum_{i=1}^N \langle b_i, \mathbf{x} \rangle \sigma'(\langle \mathbf{w}_i^0, \mathbf{x} \rangle)}_{\text{Neural Tangent (NT) model}} \end{aligned}$$

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

Lazy Training Regime

Two-layers Neural Network (NN):

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

Lazy Training regime: for some initialization $\theta_i^0 = (a_i^0, \mathbf{w}_i^0)$

$$\begin{aligned} f_{\text{NN}}(\mathbf{x}; \theta^t) &\approx f_{\text{NN}}(\mathbf{x}; \theta^0) + \langle \theta^t - \theta^0, \nabla_{\theta} f_{\text{NN}}(\mathbf{x}; \theta^0) \rangle \\ &\approx 0 + \underbrace{\sum_{i=1}^N t_i \sigma(\langle \mathbf{w}_i^0, \mathbf{x} \rangle)}_{\text{Random Features (RF) model}} + \underbrace{\sum_{i=1}^N \langle b_i, \mathbf{x} \rangle \sigma'(\langle \mathbf{w}_i^0, \mathbf{x} \rangle)}_{\text{Neural Tangent (NT) model}} \end{aligned}$$

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

Lazy Training Regime

Two-layers Neural Network (NN):

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

Lazy Training regime: for some initialization $\theta_i^0 = (a_i^0, \mathbf{w}_i^0)$

$$\begin{aligned} f_{\text{NN}}(\mathbf{x}; \theta^t) &\approx f_{\text{NN}}(\mathbf{x}; \theta^0) + \langle \theta^t - \theta^0, \nabla_{\theta} f_{\text{NN}}(\mathbf{x}; \theta^0) \rangle \\ &\approx 0 + \underbrace{\sum_{i=1}^N t_i \sigma(\langle \mathbf{w}_i^0, \mathbf{x} \rangle)}_{\text{Random Features (RF) model}} + \underbrace{\sum_{i=1}^N \langle b_i, \mathbf{x} \rangle \sigma'(\langle \mathbf{w}_i^0, \mathbf{x} \rangle)}_{\text{Neural Tangent (NT) model}} \end{aligned}$$

[Jacot, Gabriel, Hongler, 2018; Du, Zhai, Póczos, Singh, 2018; Chizat, Bach, 2018b; Arora, Du, Hu, Li, Wang, 2019; Allen-Zhu, Li, Song, 2018; Yehudai, Shamir, 2019; ...]

Lazy Training Regime

Two-layers Neural Network (NN):

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

Lazy Training regime: for some initialization $\theta_i^0 = (a_i^0, \mathbf{w}_i^0)$

$$\begin{aligned} f_{\text{NN}}(\mathbf{x}; \theta^t) &\approx f_{\text{NN}}(\mathbf{x}; \theta^0) + \langle \theta^t - \theta^0, \nabla_{\theta} f_{\text{NN}}(\mathbf{x}; \theta^0) \rangle \\ &\approx 0 + \underbrace{\sum_{i=1}^N t_i \sigma(\langle \mathbf{w}_i^0, \mathbf{x} \rangle)}_{\text{Random Features (RF) model}} + \underbrace{\sum_{i=1}^N \langle b_i, \mathbf{x} \rangle \sigma'(\langle \mathbf{w}_i^0, \mathbf{x} \rangle)}_{\text{Neural Tangent (NT) model}} \end{aligned}$$

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

Questions:

- ▶ Do RF/NT provide a good approximation to effectively trained NN?
- ▶ Do RF/NT learn effective/smart representations of the data?

Setting:

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

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

- ▶ Here $\sigma(x) = x^2$ (cf. paper for generalization)
- ▶ The neural network NN is trained by SGD
- ▶ Compare population squared error loss

Questions:

- ▶ Do RF/NT provide a good approximation to effectively trained NN?
- ▶ Do RF/NT learn effective/smart representations of the data?

Setting:

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

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

- ▶ Here $\sigma(x) = x^2$ (cf. paper for generalization)
- ▶ The neural network NN is trained by SGD
- ▶ Compare population squared error loss

Questions:

- ▶ Do RF/NT provide a good approximation to effectively trained NN?
- ▶ Do RF/NT learn effective/smart representations of the data?

Setting:

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

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

- ▶ Here $\sigma(x) = x^2$ (cf. paper for generalization)
- ▶ The neural network NN is trained by SGD
- ▶ Compare population squared error loss

Questions:

- ▶ Do RF/NT provide a good approximation to effectively trained NN?
- ▶ Do RF/NT learn effective/smart representations of the data?

Setting:

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

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

- ▶ Here $\sigma(x) = x^2$ (cf. paper for generalization)
- ▶ The neural network NN is trained by SGD
- ▶ Compare population squared error loss

Questions:

- ▶ Do RF/NT provide a good approximation to effectively trained NN?
- ▶ Do RF/NT learn effective/smart representations of the data?

Setting:

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

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

- ▶ Here $\sigma(x) = x^2$ (cf. paper for generalization)
- ▶ The neural network NN is trained by SGD
- ▶ Compare population squared error loss

Questions:

- ▶ Do RF/NT provide a good approximation to effectively trained NN?
- ▶ Do RF/NT learn effective/smart representations of the data?

Setting:

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

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

- ▶ Here $\sigma(x) = x^2$ (cf. paper for generalization)
- ▶ The neural network NN is trained by SGD
- ▶ Compare population squared error loss

Questions:

- ▶ Do RF/NT provide a good approximation to effectively trained NN?
- ▶ Do RF/NT learn effective/smart representations of the data?

Setting:

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

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

- ▶ Here $\sigma(x) = x^2$ (cf. paper for generalization)
- ▶ The neural network NN is trained by SGD
- ▶ Compare population squared error loss

Questions:

- ▶ Do RF/NT provide a good approximation to effectively trained NN?
- ▶ Do RF/NT learn effective/smart representations of the data?

Setting:

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

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

- ▶ Here $\sigma(x) = x^2$ (cf. paper for generalization)
- ▶ The neural network NN is trained by SGD
- ▶ Compare population squared error loss

Results

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

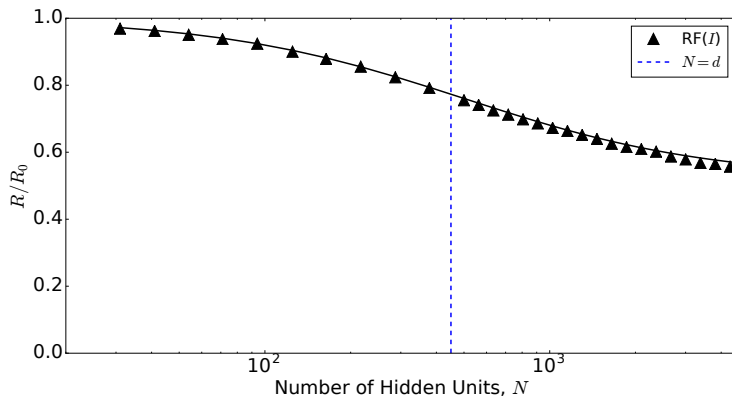


Figure: Lines are analytical predictions and dots are empirical results.

Results

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

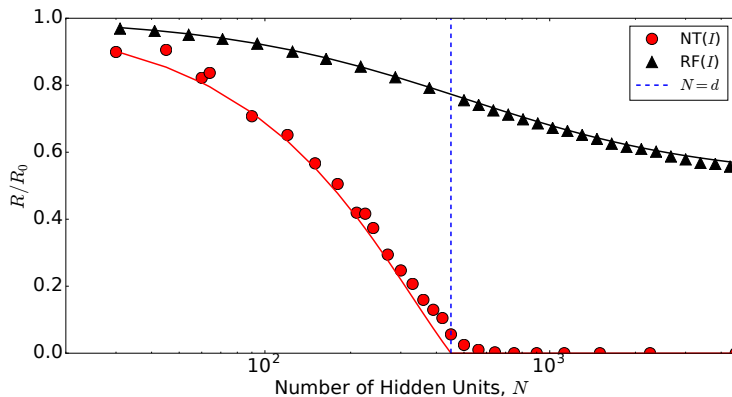


Figure: Lines are analytical predictions and dots are empirical results.

Results

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

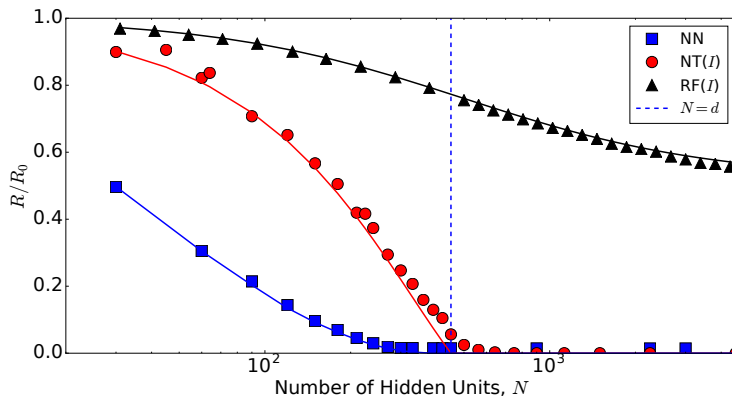


Figure: Lines are analytical predictions and dots are empirical results.

Interpretation

- ▶ RF model does not capture quadratic functions (regardless of the non-linearity)
- ▶ The NT model fits random directions spanned by $(\mathbf{w}_1^0, \dots, \mathbf{w}_N^0)$
- ▶ Fully trained NN learns the most important eigendirections
- ▶ $\exists B$ arbitrarily large gap between NN and NT

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

These phenomena are more general: mixture of Gaussians, ReLu activation...

Interpretation

- ▶ RF model does not capture quadratic functions (regardless of the non-linearity)
- ▶ The NT model fits random directions spanned by $(\mathbf{w}_1^0, \dots, \mathbf{w}_N^0)$
- ▶ Fully trained NN learns the most important eigendirections
- ▶ $\exists B$ arbitrarily large gap between NN and NT

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

These phenomena are more general: mixture of Gaussians, ReLU activation...

Interpretation

- ▶ RF model does not capture quadratic functions (regardless of the non-linearity)
- ▶ The NT model fits random directions spanned by $(\mathbf{w}_1^0, \dots, \mathbf{w}_N^0)$
- ▶ Fully trained NN learns the most important eigendirections
- ▶ $\exists B$ arbitrarily large gap between NN and NT

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

These phenomena are more general: mixture of Gaussians, ReLu activation...

Interpretation

- ▶ RF model does not capture quadratic functions (regardless of the non-linearity)
- ▶ The NT model fits random directions spanned by $(\mathbf{w}_1^0, \dots, \mathbf{w}_N^0)$
- ▶ Fully trained NN learns the most important eigendirections
- ▶ $\exists B$ arbitrarily large gap between NN and NT

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

These phenomena are more general: mixture of Gaussians, ReLu activation...

Interpretation

- ▶ RF model does not capture quadratic functions (regardless of the non-linearity)
- ▶ The NT model fits random directions spanned by $(\mathbf{w}_1^0, \dots, \mathbf{w}_N^0)$
- ▶ Fully trained NN learns the most important eigendirections
- ▶ $\exists B$ arbitrarily large gap between NN and NT

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

These phenomena are more general: mixture of Gaussians, ReLu activation...

Interpretation

- ▶ RF model does not capture quadratic functions (regardless of the non-linearity)
- ▶ The NT model fits random directions spanned by $(\mathbf{w}_1^0, \dots, \mathbf{w}_N^0)$
- ▶ Fully trained NN learns the most important eigendirections
- ▶ $\exists B$ arbitrarily large gap between NN and NT

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

These phenomena are more general: mixture of Gaussians, ReLU activation...

Interpretation

- ▶ RF model does not capture quadratic functions (regardless of the non-linearity)
- ▶ The NT model fits random directions spanned by $(\mathbf{w}_1^0, \dots, \mathbf{w}_N^0)$
- ▶ Fully trained NN learns the most important eigendirections
- ▶ $\exists B$ arbitrarily large gap between NN and NT

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

These phenomena are 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

‘Limitations of Lazy Training of Two-layers Neural Networks’

Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, Andrea Montanari