

Analyzing Finite Neural Networks: Can We Trust Neural Tangent Kernel Theory?

Mariia Seleznova & Gitta Kutyniok
(Ludwig-Maximilians-Universität München)

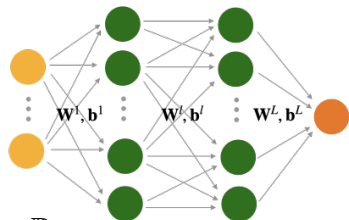
MSML 2021
August 16-19, 2021



Deep Neural Networks (DNNs)

Definition: Assume the following notation:

- ▶ Number of layers $L \geq 2$.
- ▶ Layers' widths M_ℓ , $\ell = 0, \dots, L$.
- ▶ *Weights* $W^\ell \in \mathbb{R}^{M_\ell \times M_{\ell-1}}$, $\ell \geq 1$.
- ▶ *Biases* $b^\ell \in \mathbb{R}^{M_\ell}$, $\ell \geq 1$.
- ▶ (Non-linear) *activation function* $\phi : \mathbb{R} \rightarrow \mathbb{R}$.



Then a *deep neural network (DNN)* is a function $f : \mathbb{R}^{M_0} \rightarrow \mathbb{R}^{M^L}$:

$$f(x) = W^L \phi(W^{L-1} \phi(W^{L-2} \phi(W^{L-3} \dots) + b^{L-1}) + b^L.$$

DNNs' training dynamics

Consider training a DNN with parameters $\theta = \{(W^\ell, b^\ell)\}_{\ell=1,\dots,L}$ on dataset $D = (X, Y)$, $X \in \mathbb{R}^{N \times M_0}$, $Y \in \mathbb{R}^{N \times M_L}$ by *gradient flow* in time t with *loss function* \mathcal{L} :

$$\dot{\theta}^{(t)} = -\nabla_{\theta} \mathcal{L}(f^{(t)}(X), Y)$$

DNNs' training dynamics

Consider training a DNN with parameters $\theta = \{(W^\ell, b^\ell)\}_{\ell=1, \dots, L}$ on dataset $D = (X, Y)$, $X \in \mathbb{R}^{N \times M_0}$, $Y \in \mathbb{R}^{N \times M_L}$ by *gradient flow* in time t with *loss function* \mathcal{L} :

$$\dot{\theta}^{(t)} = -\nabla_{\theta} \mathcal{L}(f^{(t)}(X), Y)$$

Then the dynamics of the output function on any input $\tilde{X} \in \mathbb{R}^{\tilde{N} \times M_0}$ is given by:

$$\dot{f}^{(t)}(\tilde{X}) = \nabla_{\theta} f^{(t)}(\tilde{X}) \cdot \dot{\theta}^{(t)}$$

DNNs' training dynamics

Consider training a DNN with parameters $\theta = \{(W^\ell, b^\ell)\}_{\ell=1,\dots,L}$ on dataset $D = (X, Y)$, $X \in \mathbb{R}^{N \times M_0}$, $Y \in \mathbb{R}^{N \times M_L}$ by *gradient flow* in time t with *loss function* \mathcal{L} :

$$\dot{\theta}^{(t)} = -\nabla_{\theta} \mathcal{L}(f^{(t)}(X), Y)$$

Then the dynamics of the output function on any input $\tilde{X} \in \mathbb{R}^{\tilde{N} \times M_0}$ is given by:

$$\dot{f}^{(t)}(\tilde{X}) = \nabla_{\theta} f^{(t)}(\tilde{X}) \cdot \dot{\theta}^{(t)}$$

Challenges:

- ▶ No analytical solutions for $f^{(t)}$ in general.

DNNs' training dynamics

Consider training a DNN with parameters $\theta = \{(W^\ell, b^\ell)\}_{\ell=1,\dots,L}$ on dataset $D = (X, Y)$, $X \in \mathbb{R}^{N \times M_0}$, $Y \in \mathbb{R}^{N \times M_L}$ by *gradient flow* in time t with *loss function* \mathcal{L} :

$$\dot{\theta}^{(t)} = -\nabla_{\theta} \mathcal{L}(f^{(t)}(X), Y)$$

Then the dynamics of the output function on any input $\tilde{X} \in \mathbb{R}^{\tilde{N} \times M_0}$ is given by:

$$\dot{f}^{(t)}(\tilde{X}) = \nabla_{\theta} f^{(t)}(\tilde{X}) \cdot \dot{\theta}^{(t)}$$

Challenges:

- ▶ No analytical solutions for $f^{(t)}$ in general.
- ▶ No access to *generalization error* $\mathbb{E}_{x,y}[\mathcal{L}(f^{(t)}(x), y)]$.

DNNs' training dynamics

Consider training a DNN with parameters $\theta = \{(W^\ell, b^\ell)\}_{\ell=1,\dots,L}$ on dataset $D = (X, Y)$, $X \in \mathbb{R}^{N \times M_0}$, $Y \in \mathbb{R}^{N \times M_L}$ by *gradient flow* in time t with *loss function* \mathcal{L} :

$$\dot{\theta}^{(t)} = -\nabla_{\theta} \mathcal{L}(f^{(t)}(X), Y)$$

Then the dynamics of the output function on any input $\tilde{X} \in \mathbb{R}^{\tilde{N} \times M_0}$ is given by:

$$\dot{f}^{(t)}(\tilde{X}) = \nabla_{\theta} f^{(t)}(\tilde{X}) \cdot \dot{\theta}^{(t)}$$

Challenges:

- ▶ No analytical solutions for $f^{(t)}$ in general.
- ▶ No access to *generalization error* $\mathbb{E}_{x,y}[\mathcal{L}(f^{(t)}(x), y)]$.
- ▶ No access to model's stability and robustness.

DNNs' training dynamics

Consider training a DNN with parameters $\theta = \{(W^\ell, b^\ell)\}_{\ell=1, \dots, L}$ on dataset $D = (X, Y)$, $X \in \mathbb{R}^{N \times M_0}$, $Y \in \mathbb{R}^{N \times M_L}$ by *gradient flow* in time t with *loss function* \mathcal{L} :

$$\dot{\theta}^{(t)} = -\nabla_{\theta} \mathcal{L}(f^{(t)}(X), Y)$$

Then the dynamics of the output function on any input $\tilde{X} \in \mathbb{R}^{\tilde{N} \times M_0}$ is given by:

$$\dot{f}^{(t)}(\tilde{X}) = \nabla_{\theta} f^{(t)}(\tilde{X}) \cdot \dot{\theta}^{(t)}$$

Challenges:

- ▶ No analytical solutions for $f^{(t)}$ in general.
- ▶ No access to *generalization error* $\mathbb{E}_{x,y}[\mathcal{L}(f^{(t)}(x), y)]$.
- ▶ No access to model's stability and robustness.

↪ *Neural Tangent Kernel (NTK) theory addresses these challenges in a special case of infinitely-wide DNNs!*

Neural Tangent Kernel Theory

Consider *squared loss* $\mathcal{L}(\hat{Y}, Y) = \frac{1}{2N} \|(\hat{Y} - Y)\|_2^2$ and for simplicity set $M_L = 1$. Then the gradient flow dynamics of a DNN takes form:

$$\dot{\theta}^{(t)} = -\nabla_{\theta} \mathcal{L}(f^{(t)}(X), Y) = -\frac{1}{N} \nabla_{\theta} f^{(t)}(X)^T \cdot (f^{(t)}(X) - Y),$$

$$\dot{f}^{(t)}(\tilde{X}) = \nabla_{\theta} f^{(t)}(\tilde{X}) \cdot \dot{\theta}^{(t)} = -\frac{1}{N} \underbrace{\nabla_{\theta} f^{(t)}(\tilde{X}) \nabla_{\theta} f^{(t)}(X)^T}_{\Theta^{(t)}(\tilde{X}, X)} \cdot (f^{(t)}(X) - Y).$$

Neural Tangent Kernel Theory

Consider *squared loss* $\mathcal{L}(\hat{Y}, Y) = \frac{1}{2N} \|(\hat{Y} - Y)\|_2^2$ and for simplicity set $M_L = 1$. Then the gradient flow dynamics of a DNN takes form:

$$\dot{\theta}^{(t)} = -\nabla_{\theta} \mathcal{L}(f^{(t)}(X), Y) = -\frac{1}{N} \nabla_{\theta} f^{(t)}(X)^T \cdot (f^{(t)}(X) - Y),$$

$$\dot{f}^{(t)}(\tilde{X}) = \nabla_{\theta} f^{(t)}(\tilde{X}) \cdot \dot{\theta}^{(t)} = -\frac{1}{N} \underbrace{\nabla_{\theta} f^{(t)}(\tilde{X}) \nabla_{\theta} f^{(t)}(X)^T}_{\Theta^{(t)}(\tilde{X}, X)} \cdot (f^{(t)}(X) - Y).$$

Definition: *Neural tangent kernel (NTK)* of a DNN with output function f and trainable parameters θ is given by

$$\Theta(x_i, x_j) := \nabla_{\theta} f(x_i)^T \nabla_{\theta} f(x_j), \quad x_i, x_j \in \mathbb{R}^{M_0}.$$

Neural Tangent Kernel Theory

Results on infinite-width limit of NTK $M_\ell \rightarrow \infty, \ell = 1, \dots, L - 1$:^[1]

Neural Tangent Kernel Theory

Results on infinite-width limit of NTK $M_\ell \rightarrow \infty, \ell = 1, \dots, L - 1$:^[1]

- ▶ NTK is *deterministic under random initialization*:

$$\Theta^{(0)}(x_i, x_j) \rightarrow \mathbb{E}_\theta[\Theta^{(0)}(x_i, x_j)] = \Theta^*(x_i, x_j),$$

$$\text{where } W_{ij}^\ell = \frac{\sigma_w}{\sqrt{M^\ell}} w_{ij}^\ell, \quad w_{ij}^\ell \sim \mathcal{N}(0, 1),$$

$$b_i^\ell = \sigma_b \beta_i^\ell, \quad \beta_i^\ell \sim \mathcal{N}(0, 1).$$

Neural Tangent Kernel Theory

Results on infinite-width limit of NTK $M_\ell \rightarrow \infty, \ell = 1, \dots, L - 1$:^[1]

- ▶ NTK is *deterministic under random initialization*:

$$\Theta^{(0)}(x_i, x_j) \rightarrow \mathbb{E}_\theta[\Theta^{(0)}(x_i, x_j)] = \Theta^*(x_i, x_j),$$

$$\text{where } W_{ij}^\ell = \frac{\sigma_w}{\sqrt{M^\ell}} w_{ij}^\ell, \quad w_{ij}^\ell \sim \mathcal{N}(0, 1),$$

$$b_i^\ell = \sigma_b \beta_i^\ell, \quad \beta_i^\ell \sim \mathcal{N}(0, 1).$$

- ▶ NTK stays *constant during training*:

$$\Theta^{(t)}(x_i, x_j) \rightarrow \Theta^*(x_i, x_j).$$

Neural Tangent Kernel Theory

Results on infinite-width limit of NTK $M_\ell \rightarrow \infty, \ell = 1, \dots, L - 1$:^[1]

- ▶ NTK is *deterministic under random initialization*:

$$\Theta^{(0)}(x_i, x_j) \rightarrow \mathbb{E}_\theta[\Theta^{(0)}(x_i, x_j)] = \Theta^*(x_i, x_j),$$

$$\text{where } W_{ij}^\ell = \frac{\sigma_w}{\sqrt{M^\ell}} w_{ij}^\ell, \quad w_{ij}^\ell \sim \mathcal{N}(0, 1),$$

$$b_i^\ell = \sigma_b \beta_i^\ell, \quad \beta_i^\ell \sim \mathcal{N}(0, 1).$$

- ▶ NTK stays *constant during training*:

$$\Theta^{(t)}(x_i, x_j) \rightarrow \Theta^*(x_i, x_j).$$

Then in the infinite-width limit gradient flow dynamics with squared loss:

$$\dot{f}^{(t)}(\tilde{X}) = -\frac{1}{N} \Theta^*(\tilde{X}, X) \cdot (f^{(t)}(X) - Y)$$

Neural Tangent Kernel Theory

Results on infinite-width limit of NTK $M_\ell \rightarrow \infty, \ell = 1, \dots, L - 1$:^[1]

- ▶ NTK is *deterministic under random initialization*:

$$\Theta^{(0)}(x_i, x_j) \rightarrow \mathbb{E}_\theta[\Theta^{(0)}(x_i, x_j)] = \Theta^*(x_i, x_j),$$

$$\text{where } W_{ij}^\ell = \frac{\sigma_w}{\sqrt{M^\ell}} w_{ij}^\ell, \quad w_{ij}^\ell \sim \mathcal{N}(0, 1),$$

$$b_i^\ell = \sigma_b \beta_i^\ell, \quad \beta_i^\ell \sim \mathcal{N}(0, 1).$$

- ▶ NTK stays *constant during training*:

$$\Theta^{(t)}(x_i, x_j) \rightarrow \Theta^*(x_i, x_j).$$

Then in the infinite-width limit gradient flow dynamics with squared loss:

$$\dot{f}^{(t)}(\tilde{X}) = -\frac{1}{N} \Theta^*(\tilde{X}, X) \cdot (f^{(t)}(X) - Y)$$

\rightsquigarrow *Infinitely-wide DNNs evolve as kernel regression with NTK kernel!*

Do finite DNNs behave as infinite-width ones?

Problems:

- ▶ If NTK matrix is constant, *no feature learning occurs*.

Do finite DNNs behave as infinite-width ones?

Problems:

- ▶ If NTK matrix is constant, *no feature learning occurs*.
- ▶ Empirical performance of NTK and finite DNNs differs.^{[2],[3]}

Do finite DNNs behave as infinite-width ones?

Problems:

- ▶ If NTK matrix is constant, *no feature learning occurs*.
- ▶ Empirical performance of NTK and finite DNNs differs.^{[2],[3]}
- ▶ In infinite-*depth-and-width* limit ($L/M > 0$), NTK of ReLU DNNs initialized with $\sigma_w^2 = 2, \sigma_b^2 = 0$ is random.^[4]

Do finite DNNs behave as infinite-width ones?

Problems:

- ▶ If NTK matrix is constant, *no feature learning occurs*.
- ▶ Empirical performance of NTK and finite DNNs differs.^{[2],[3]}
- ▶ In infinite-*depth-and-width* limit ($L/M > 0$), NTK of ReLU DNNs initialized with $\sigma_w^2 = 2, \sigma_b^2 = 0$ is random.^[4]

~> *It is not clear when NTK theory explains DNNs' behavior!*

Do finite DNNs behave as infinite-width ones?

Problems:

- ▶ If NTK matrix is constant, *no feature learning occurs*.
- ▶ Empirical performance of NTK and finite DNNs differs.^{[2],[3]}
- ▶ In infinite-*depth-and-width* limit ($L/M > 0$), NTK of ReLU DNNs initialized with $\sigma_w^2 = 2, \sigma_b^2 = 0$ is random.^[4]

~> *It is not clear when NTK theory explains DNNs' behavior!*

Our contributions:

- ▶ Study *ReLU* and *sigmoid* DNNs with various hyperparameters (σ_w, σ_b, L, M).

Do finite DNNs behave as infinite-width ones?

Problems:

- ▶ If NTK matrix is constant, *no feature learning occurs*.
- ▶ Empirical performance of NTK and finite DNNs differs.^{[2],[3]}
- ▶ In infinite-*depth-and-width* limit ($L/M > 0$), NTK of ReLU DNNs initialized with $\sigma_w^2 = 2, \sigma_b^2 = 0$ is random.^[4]

↪ *It is not clear when NTK theory explains DNNs' behavior!*

Our contributions:

- ▶ Study *ReLU* and *sigmoid* DNNs with various hyperparameters (σ_w, σ_b, L, M).
- ▶ Identify *two phases* in hyperparameter space where NTK regime does and does not hold.

Do finite DNNs behave as infinite-width ones?

Problems:

- ▶ If NTK matrix is constant, *no feature learning occurs*.
- ▶ Empirical performance of NTK and finite DNNs differs.^{[2],[3]}
- ▶ In infinite-*depth-and-width* limit ($L/M > 0$), NTK of ReLU DNNs initialized with $\sigma_w^2 = 2, \sigma_b^2 = 0$ is random.^[4]

↪ *It is not clear when NTK theory explains DNNs' behavior!*

Our contributions:

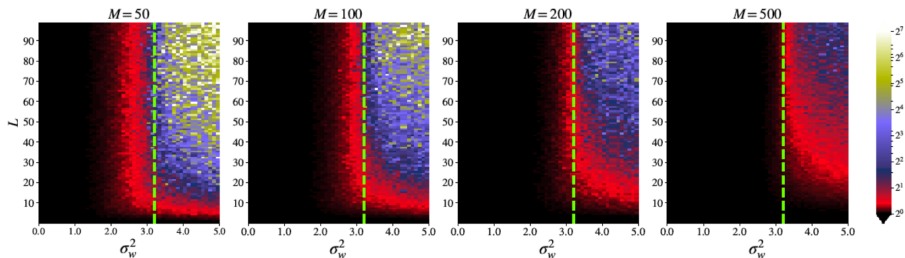
- ▶ Study *ReLU* and *sigmoid* DNNs with various hyperparameters (σ_w, σ_b, L, M).
- ▶ Identify *two phases* in hyperparameter space where NTK regime does and does not hold.
- ▶ Study variance of DNNs output $\text{Var}_{\theta, D} [f^{(t \rightarrow \infty)}(x)]$ under NTK theory.

Randomness at initialization

Setup:

- ▶ Fully-connected *tanh* networks with L layers and constant width M .
- ▶ Initialized as $W_{ij}^\ell \sim \mathcal{N}(0, \frac{\sigma_w^2}{M})$, $b_i^\ell \sim \mathcal{N}(0, \sigma_b^2)$

$\frac{\mathbb{E}_\theta[\Theta^{(0)}(x,x)^2]}{\mathbb{E}_\theta^2[\Theta^{(0)}(x,x)]}$ ratio measures randomness at initialization:



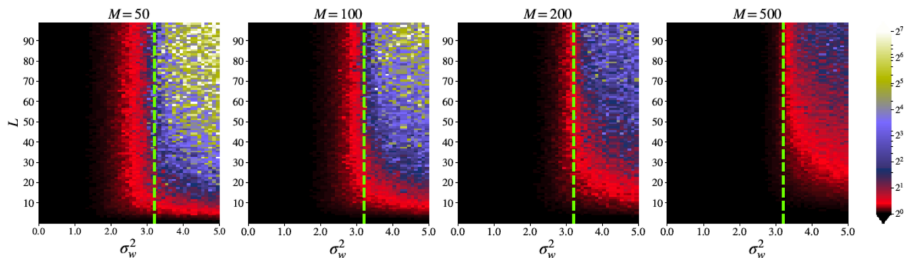
We use standard parametrization instead of NTK parametrization here. However, for constant-width networks this does not affect the results.

Randomness at initialization

Setup:

- ▶ Fully-connected *tanh* networks with L layers and constant width M .
- ▶ Initialized as $W_{ij}^\ell \sim \mathcal{N}(0, \frac{\sigma_w^2}{M})$, $b_i^\ell \sim \mathcal{N}(0, \sigma_b^2)$

$\frac{\mathbb{E}_\theta[\Theta^{(0)}(x,x)^2]}{\mathbb{E}_\theta^2[\Theta^{(0)}(x,x)]}$ ratio measures randomness at initialization:



\leadsto Deep NNs with large σ_w are random at initialization!

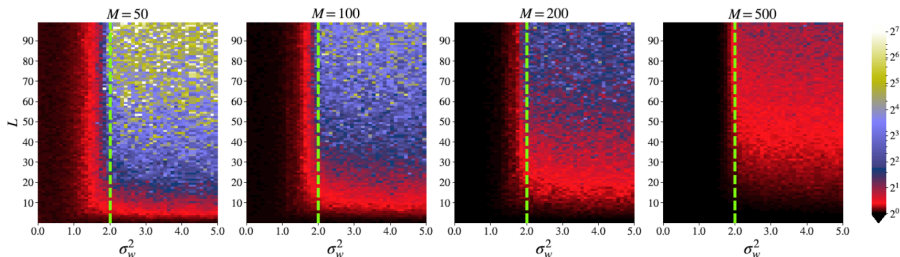
We use standard parametrization instead of NTK parametrization here. However, for constant-width networks this does not affect the results.

Randomness at initialization

Setup:

- ▶ Fully-connected *ReLU* networks with L layers and constant width M .
- ▶ Initialized as $W_{ij}^\ell \sim \mathcal{N}(0, \frac{\sigma_w^2}{M})$, $b_i^\ell \sim \mathcal{N}(0, \sigma_b^2)$

$\frac{\mathbb{E}_\theta[\Theta^{(0)}(x,x)^2]}{\mathbb{E}_\theta^2[\Theta^{(0)}(x,x)]}$ ratio measures randomness at initialization:



↷ *Deep NNs with large σ_w are random at initialization!*

We use standard parametrization instead of NTK parametrization here. However, for constant-width networks this does not affect the results.

Vanishing and exploding gradients

Behaviour of *gradients at initialization* is controlled by variable χ :^[5]

$$\chi := \sigma_w^2 \int [\phi'(\sqrt{q^*}v)]^2 Dv, \quad Dv = \frac{dv}{\sqrt{2\pi}} e^{-v^2/2},$$

where $q^* = \lim_{\ell \rightarrow \infty} q^\ell$ and $q^\ell := \frac{1}{M_\ell} \sum_{k=1}^{M_\ell} (z_k^\ell)^2$ is the pre-activation “length” in layer ℓ .

Vanishing and exploding gradients

Behaviour of *gradients at initialization* is controlled by variable χ :^[5]

$$\chi := \sigma_w^2 \int [\phi'(\sqrt{q^*}v)]^2 Dv, \quad Dv = \frac{dv}{\sqrt{2\pi}} e^{-v^2/2},$$

where $q^* = \lim_{\ell \rightarrow \infty} q^\ell$ and $q^\ell := \frac{1}{M_\ell} \sum_{k=1}^{M_\ell} (z_k^\ell)^2$ is the pre-activation “length” in layer ℓ .

$\leadsto \chi$ depends on hyperparameters (σ_w, σ_b) and activation function ϕ .

Vanishing and exploding gradients

Behaviour of *gradients at initialization* is controlled by variable χ :^[5]

$$\chi := \sigma_w^2 \int [\phi'(\sqrt{q^*}v)]^2 Dv, \quad Dv = \frac{dv}{\sqrt{2\pi}} e^{-v^2/2},$$

where $q^* = \lim_{\ell \rightarrow \infty} q^\ell$ and $q^\ell := \frac{1}{M_\ell} \sum_{k=1}^{M_\ell} (z_k^\ell)^2$ is the pre-activation “length” in layer ℓ .

$\leadsto \chi$ depends on hyperparameters (σ_w, σ_b) and activation function ϕ .

We can identify the following situations based on χ :

- ▶ *Chaotic phase:* If $\chi > 1$, gradients explode as they backpropagate.

Vanishing and exploding gradients

Behaviour of *gradients at initialization* is controlled by variable χ :^[5]

$$\chi := \sigma_w^2 \int [\phi'(\sqrt{q^*}v)]^2 Dv, \quad Dv = \frac{dv}{\sqrt{2\pi}} e^{-v^2/2},$$

where $q^* = \lim_{\ell \rightarrow \infty} q^\ell$ and $q^\ell := \frac{1}{M_\ell} \sum_{k=1}^{M_\ell} (z_k^\ell)^2$ is the pre-activation “length” in layer ℓ .

$\leadsto \chi$ depends on hyperparameters (σ_w, σ_b) and activation function ϕ .

We can identify the following situations based on χ :

- ▶ *Chaotic phase*: If $\chi > 1$, gradients explode as they backpropagate.
- ▶ *Ordered phase*: If $\chi < 1$, gradients vanish.

$q^\ell(x)$ depends only on the norm of x . Therefore, for simplicity of notation we can assume normalized inputs and omit argument x here.

Vanishing and exploding gradients

Behaviour of *gradients at initialization* is controlled by variable χ :^[5]

$$\chi := \sigma_w^2 \int [\phi'(\sqrt{q^*}v)]^2 Dv, \quad Dv = \frac{dv}{\sqrt{2\pi}} e^{-v^2/2},$$

where $q^* = \lim_{\ell \rightarrow \infty} q^\ell$ and $q^\ell := \frac{1}{M_\ell} \sum_{k=1}^{M_\ell} (z_k^\ell)^2$ is the pre-activation “length” in layer ℓ .

$\leadsto \chi$ depends on hyperparameters (σ_w, σ_b) and activation function ϕ .

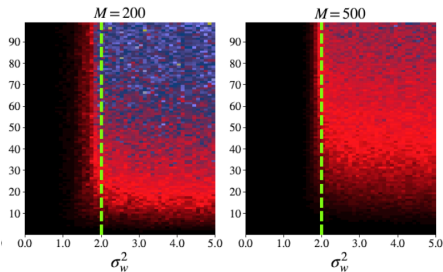
We can identify the following situations based on χ :

- ▶ *Chaotic phase*: If $\chi > 1$, gradients explode as they backpropagate.
- ▶ *Ordered phase*: If $\chi < 1$, gradients vanish.
- ▶ *«Edge of chaos» (EOC)*: $\chi \approx 1$ allows deeper signal propagation.

$q^\ell(x)$ depends only on the norm of x . Therefore, for simplicity of notation we can assume normalized inputs and omit argument x here.

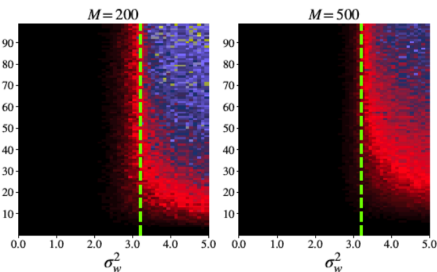
Randomness at initialization

ReLU DNNs



$$\chi = 1 \text{ if } \sigma_w^2 = 2$$

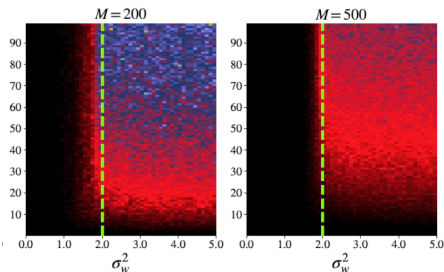
tanh DNNs



$$\chi = 1 \text{ if } \sigma_w^2 \approx 3.2$$

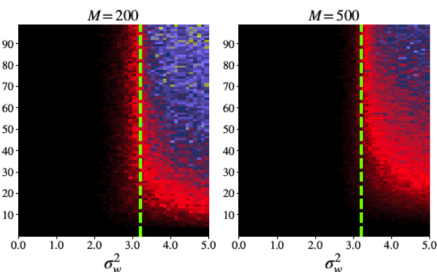
Randomness at initialization

ReLU DNNs



$$\chi = 1 \text{ if } \sigma_w^2 = 2$$

tanh DNNs

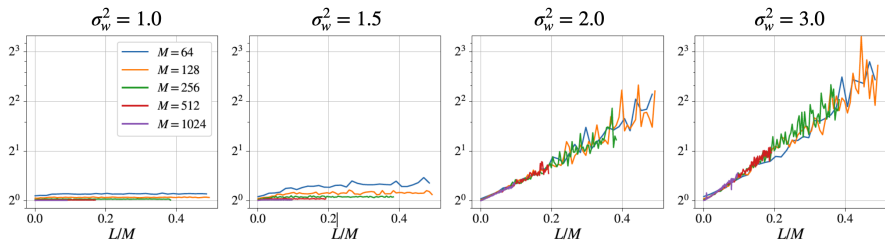


$$\chi = 1 \text{ if } \sigma_w^2 \approx 3.2$$

\leadsto Deep networks in chaotic phase are random at initialization!

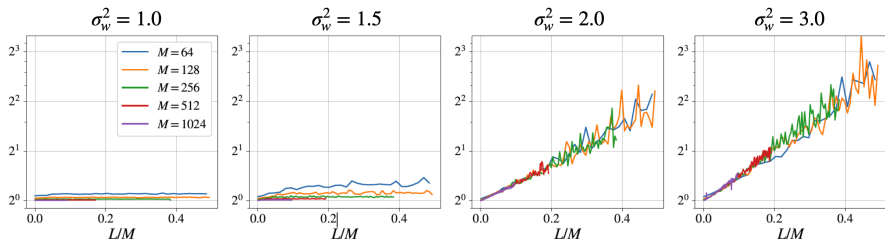
Randomness at initialization

$\frac{\mathbb{E}_\theta[\Theta^{(0)}(x,x)^2]}{\mathbb{E}_\theta^2[\Theta^{(0)}(x,x)]}$ ratio as a function of $\frac{L}{M}$:



Randomness at initialization

$\frac{\mathbb{E}_\theta [\Theta^{(0)}(x,x)^2]}{\mathbb{E}_\theta^2 [\Theta^{(0)}(x,x)]}$ ratio as a function of $\frac{L}{M}$:



\leadsto Exponential growth in L/M in the chaotic phase.

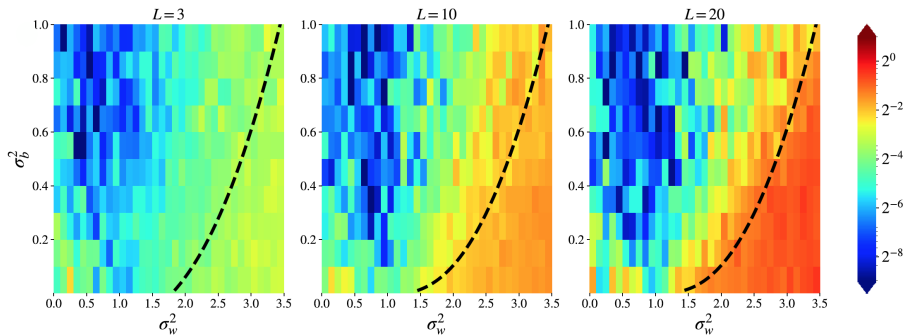
\leadsto Dependence on $1/M$ in the ordered phase.

Change during training

Setup:

- Fully-connected *tanh* networks with L layers and constant width $M = 256$.

$\frac{\|\Theta^{(t)} - \Theta^{(0)}\|_F}{\|\Theta^{(0)}\|_F}$ shows if NTK changes significantly during training:

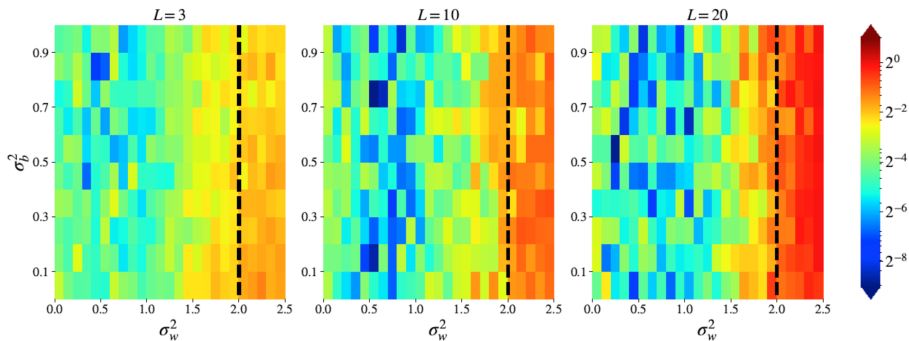


Change during training

Setup:

- Fully-connected *ReLU* networks with L layers and constant width $M = 256$.

$\frac{\|\Theta^{(t)} - \Theta^{(0)}\|_F}{\|\Theta^{(0)}\|_F}$ shows if NTK changes significantly during training:

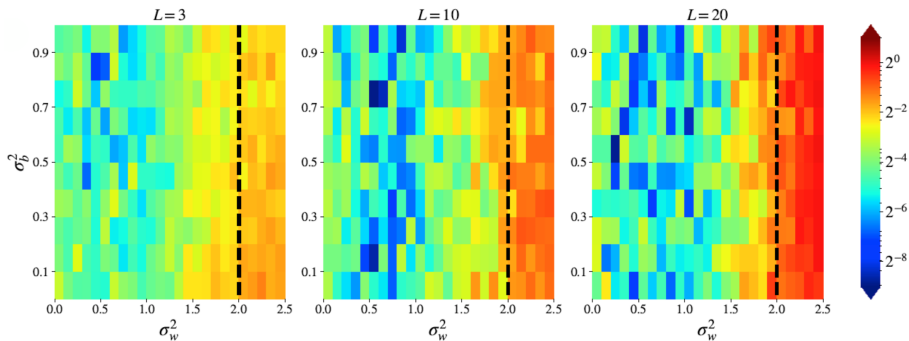


Change during training

Setup:

- Fully-connected *ReLU* networks with L layers and constant width $M = 256$.

$\frac{\|\Theta^{(t)} - \Theta^{(0)}\|_F}{\|\Theta^{(0)}\|_F}$ shows if NTK changes significantly during training:



\leadsto NTK changes significantly during training in the chaotic phase.

Generalization in the NTK regime

NTK has the following *structure* at initialization:

$$\begin{aligned}\Theta^*(X) &= \bar{\Theta}^*(\mathbb{I}_N + \epsilon(X)), \\ \bar{\Theta}^* &= (\bar{\kappa}_1 - \bar{\kappa}_2)\mathbb{I}_N + \bar{\kappa}_2\mathbb{1}_N\mathbb{1}_N^T,\end{aligned}$$

where $\epsilon(X) \xrightarrow[L \rightarrow \infty]{} 0$ [6] is the only data-dependent part and $\bar{\kappa}_i, i = 1, 2$ are controlled by depth and gradients' behaviour.

Generalization in the NTK regime

NTK has the following *structure* at initialization:

$$\Theta^*(X) = \bar{\Theta}^*(\mathbb{I}_N + \epsilon(X)),$$
$$\bar{\Theta}^* = (\bar{\kappa}_1 - \bar{\kappa}_2)\mathbb{I}_N + \bar{\kappa}_2\mathbb{1}_N\mathbb{1}_N^T,$$

where $\epsilon(X) \xrightarrow[L \rightarrow \infty]{} 0$ ^[6] is the only data-dependent part and $\bar{\kappa}_i, i = 1, 2$ are controlled by depth and gradients' behaviour.

NTK behaviour depends on initialization:

- ▶ *Chaotic phase:* $\bar{\kappa}_1/\bar{\kappa}_2 \gg 1$ for large $L \Rightarrow \Theta^* \approx \bar{\kappa}_1\mathbb{I}_N$.
- ▶ *Ordered phase:* $\bar{\kappa}_1/\bar{\kappa}_2 \approx 1$ for large $L \Rightarrow \Theta^* \approx \bar{\kappa}_2\mathbb{1}_N\mathbb{1}_N^T$.

Generalization in the NTK regime

NTK has the following *structure* at initialization:

$$\Theta^*(X) = \bar{\Theta}^*(\mathbb{I}_N + \epsilon(X)),$$
$$\bar{\Theta}^* = (\bar{\kappa}_1 - \bar{\kappa}_2)\mathbb{I}_N + \bar{\kappa}_2\mathbb{1}_N\mathbb{1}_N^T,$$

where $\epsilon(X) \xrightarrow{L \rightarrow \infty} 0$ [6] is the only data-dependent part and $\bar{\kappa}_i, i = 1, 2$ are controlled by depth and gradients' behaviour.

NTK behaviour depends on initialization:

- ▶ *Chaotic phase:* $\bar{\kappa}_1/\bar{\kappa}_2 \gg 1$ for large $L \Rightarrow \Theta^* \approx \bar{\kappa}_1\mathbb{I}_N$.
- ▶ *Ordered phase:* $\bar{\kappa}_1/\bar{\kappa}_2 \approx 1$ for large $L \Rightarrow \Theta^* \approx \bar{\kappa}_2\mathbb{1}_N\mathbb{1}_N^T$.

~> DNNs in the NTK regime have different dynamics in ordered and chaotic phases!

Generalization in the NTK regime

Theorem (Seleznova&Kutyniok, 2020): Assume the NTK matrix is well-conditioned ($\bar{\kappa}_1/\bar{\kappa}_2 \gg 1$). Then for the *variance of a trained DNN in the NTK regime* we have:

$$\text{Var}_{\theta, X}[f^{(t \rightarrow \infty)}(\tilde{x})] \approx \left(1 + \frac{A^2}{N}\right) \left(\text{Var}^{(0)} - \text{Cov}^{(0)}\right) + (A - 1)^2 \text{Cov}^{(0)},$$

where $A = \frac{N}{\bar{\kappa}_1/\bar{\kappa}_2 + (N-1)}$, $\text{Var}^{(0)} := \text{Var}_{\theta, X, \tilde{x}}[f^{(0)}(\tilde{x})]$ is the output variance at initialization, $\text{Cov}^{(0)} = \text{Cov}_{\theta, X, x_i \neq x_j}[f^{(0)}(x_i), f^{(0)}(x_j)]$ is the output covariance on two different inputs.

Generalization in the NTK regime

Theorem (Seleznova&Kutyniok, 2020): Assume the NTK matrix is well-conditioned ($\bar{\kappa}_1/\bar{\kappa}_2 \gg 1$). Then for the *variance of a trained DNN in the NTK regime* we have:

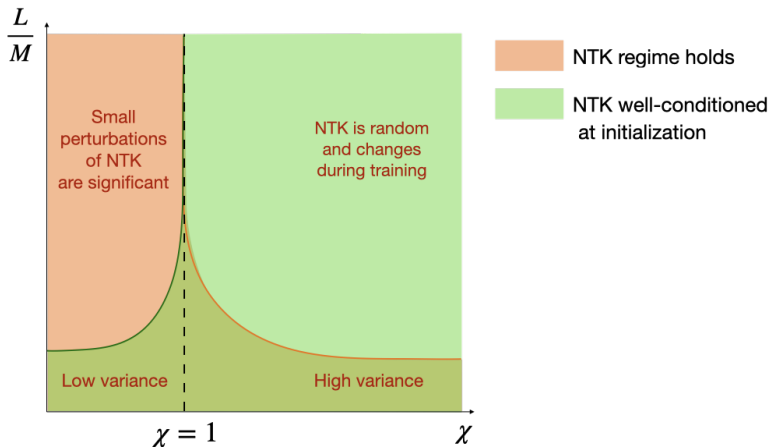
$$\text{Var}_{\theta, \mathcal{X}}[f^{(t \rightarrow \infty)}(\tilde{x})] \approx \left(1 + \frac{A^2}{N}\right) \left(\text{Var}^{(0)} - \text{Cov}^{(0)}\right) + (A - 1)^2 \text{Cov}^{(0)},$$

where $A = \frac{N}{\bar{\kappa}_1/\bar{\kappa}_2 + (N-1)}$, $\text{Var}^{(0)} := \text{Var}_{\theta, \mathcal{X}, \tilde{x}}[f^{(0)}(\tilde{x})]$ is the output variance at initialization, $\text{Cov}^{(0)} = \text{Cov}_{\theta, \mathcal{X}, x_i \neq x_j}[f^{(0)}(x_i), f^{(0)}(x_j)]$ is the output covariance on two different inputs.

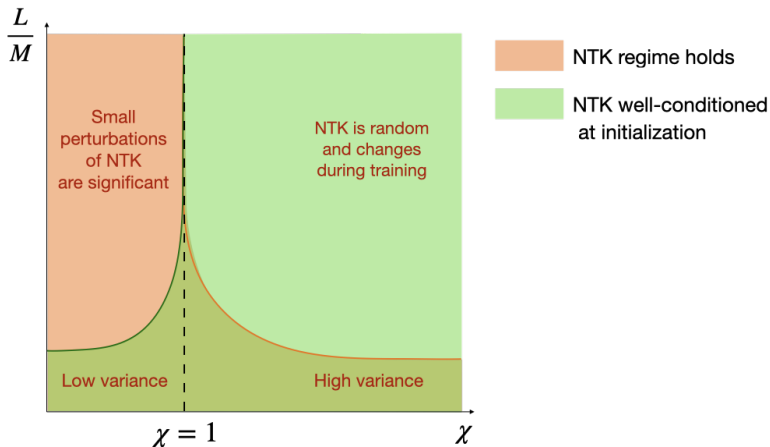
If all the conditions hold, we have:

- ▶ *Chaotic phase:* $\text{Var}_{\theta, \mathcal{X}}[f^{(t \rightarrow \infty)}(\tilde{x})] \propto \text{Var}^{(0)}$ – large variance, which growth with depth L .
- ▶ *Ordered phase:* $\text{Var}_{\theta, \mathcal{X}}[f^{(t \rightarrow \infty)}(\tilde{x})] \approx 0$ – low variance for large L .

When can we trust the results?



When can we trust the results?



↪ *Deep networks cannot be analyzed within the NTK theory!*

Conclusions

- ▶ NTK theory is a powerful tool to analyze DNNs theoretically. However, it is important to understand when it is applicable.

Conclusions

- ▶ NTK theory is a powerful tool to analyze DNNs theoretically. However, it is important to understand when it is applicable.
- ▶ *Empirical NTK* behaves as theoretical NTK for DNNs in the ordered phase but not in the chaotic phase.

Conclusions

- ▶ NTK theory is a powerful tool to analyze DNNs theoretically. However, it is important to understand when it is applicable.
- ▶ *Empirical NTK* behaves as theoretical NTK for DNNs in the ordered phase but not in the chaotic phase.
- ▶ Generalization of *shallow networks* ($L/M \approx 0$) can be analyzed within the NTK theory.

Conclusions

- ▶ NTK theory is a powerful tool to analyze DNNs theoretically. However, it is important to understand when it is applicable.
- ▶ *Empirical NTK* behaves as theoretical NTK for DNNs in the ordered phase but not in the chaotic phase.
- ▶ Generalization of *shallow networks* ($L/M \approx 0$) can be analyzed within the NTK theory.
- ▶ *Deep networks* are hard to analyze within the NTK theory.
↪ New approaches are needed to analyze DNNs theoretically.

References:

- [1] Jacot et al. *Neural Tangent Kernel: Convergence and Generalization in Neural Networks*. 2018
- [2] Lee et al. *Finite Versus Infinite Neural Networks: an Empirical Study*. 2020
- [3] Bai & Lee. *Beyond Linearization: On Quadratic and Higher-Order Approximation of Wide Neural Networks*. 2020
- [4] Hanin & Nica. *Finite Depth and Width Corrections to the Neural Tangent Kernel*. 2020
- [5] Schoenholz et al. *Deep information propagation*. 2017
- [6] Xiao et al. *Disentangling Trainability and Generalization in Deep Neural Networks*. 2020

Thank you for your attention!

Parametrization*

Infinite-width limit of NTK is normally considered in *NTK parametrization (NTP)* instead of *standard parametrization (SP)*.

$$\text{SP: } a^{l+1} = \phi\left(W^l a^l + b^l\right) \quad \text{NTP: } a^{l+1} = \phi\left(\frac{\sigma_w}{\sqrt{M^l}} w^l x^l + \sigma_b b^l\right)$$

$$W_{ij}^l \sim \mathcal{N}\left(0, \frac{\sigma_w^2}{M^l}\right), b_i^l \sim \mathcal{N}\left(0, \sigma_b^2\right) \quad w_{ij}^l \sim \mathcal{N}(0, 1), b_i^l \sim \mathcal{N}(0, 1)$$

The change from SP to NTK amounts to: $\nabla_{W^l} f^{(t)}(x) \rightarrow \frac{1}{\sqrt{M^l}} \nabla_{W^l} f^{(t)}(x)$

And for constant-width networks: $\Theta^{(t)}(x_i, x_j) \approx \frac{1}{M} \Theta^{(t)}(x_i, x_j)$

\leadsto *The same dynamics of $f^{(t)}$ with proper adjustment of η .*

\leadsto $\frac{\mathbb{E}[\Theta^{(0)}(x, x)^2]}{\mathbb{E}^2[\Theta^{(0)}(x, x)]}$ and $\frac{\|\Theta^{(t)} - \Theta^{(0)}\|_F}{\|\Theta^{(0)}\|_F}$ ratios are not affected.