

On the linearity of large non-linear models: when and why the tangent kernel is constant

Chaoyue Liu, Libin Zhu, Mikhail Belkin,

What's it about?

What's it about?

- NeurIPS paper this year

What's it about?

- NeurIPS paper this year
- Revisiting of the proof that neural tangent kernel remains **constant** during training

What's it about?

- NeurIPS paper this year
- Revisiting of the proof that neural tangent kernel remains **constant** during training
- Conceptually **easier** reasoning when and why NTK is constant

What's it about?

- NeurIPS paper this year
- Revisiting of the proof that neural tangent kernel remains **constant** during training
- Conceptually **easier** reasoning when and why NTK is constant
- **Lazy training** is not the origin of constancy, it's a consequence of the architecture!

Setup

Setup

- As always, fully-connected network of L layers mapping to \mathbb{R} :

Setup

- As always, fully-connected network of L layers mapping to \mathbb{R} :
 - $\alpha^{(0)}(\mathbf{x}) = \mathbf{x}$

Setup

- As always, fully-connected network of L layers mapping to \mathbb{R} :
 - $\alpha^{(0)}(\mathbf{x}) = \mathbf{x}$
 - $\alpha^{(l)}(\mathbf{x}) = \sigma(\tilde{\alpha}^{(l-1)}(\mathbf{x}))$

Setup

- As always, fully-connected network of L layers mapping to \mathbb{R} :
 - $\alpha^{(0)}(\mathbf{x}) = \mathbf{x}$
 - $\alpha^{(l)}(\mathbf{x}) = \sigma(\tilde{\alpha}^{(l-1)}(\mathbf{x}))$
 - $\tilde{\alpha}^{(l)}(\mathbf{x}) = \frac{1}{\sqrt{m_{l-1}}} \mathbf{W}^{(l)} \alpha^{(l-1)}(\mathbf{x})$

Setup

- As always, fully-connected network of L layers mapping to \mathbb{R} :
 - $\alpha^{(0)}(\mathbf{x}) = \mathbf{x}$
 - $\alpha^{(l)}(\mathbf{x}) = \sigma(\tilde{\alpha}^{(l-1)}(\mathbf{x}))$
 - $\tilde{\alpha}^{(l)}(\mathbf{x}) = \frac{1}{\sqrt{m_{l-1}}} \mathbf{W}^{(l)} \alpha^{(l-1)}(\mathbf{x})$

where $\mathbf{W}^{(l)} \in \mathbb{R}^{m_{l-1} \times m_l}$ and $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ some **non-linearity**

Setup

- As always, fully-connected network of L layers mapping to \mathbb{R} :
 - $\alpha^{(0)}(\mathbf{x}) = \mathbf{x}$
 - $\alpha^{(l)}(\mathbf{x}) = \sigma(\tilde{\alpha}^{(l-1)}(\mathbf{x}))$
 - $\tilde{\alpha}^{(l)}(\mathbf{x}) = \frac{1}{\sqrt{m_{l-1}}} \mathbf{W}^{(l)} \alpha^{(l-1)}(\mathbf{x})$

where $\mathbf{W}^{(l)} \in \mathbb{R}^{m_{l-1} \times m_l}$ and $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ some **non-linearity**

- The final output is given by

$$f(\mathbf{x}) = \frac{1}{\sqrt{m}} \mathbf{v}^T \alpha^{(L)}(\mathbf{x})$$

where $\mathbf{v} \in \mathbb{R}^m$ is **not** trainable

Initialization and Conditions

Initialization and Conditions

Initialize all the parameters in the network as follows:

Initialization and Conditions

Initialize all the parameters in the network as follows:

- $W_{ij}^{(l)} \overset{i.i.d.}{\sim} \mathcal{N}(0, 1)$

Initialization and Conditions

Initialize all the parameters in the network as follows:

- $W_{ij}^{(l)} \overset{i.i.d.}{\sim} \mathcal{N}(0, 1)$
- $v_i \sim \mathcal{U}(\{-1, 1\})$

Initialization and Conditions

Initialize all the parameters in the network as follows:

- $W_{ij}^{(l)} \overset{i.i.d.}{\sim} \mathcal{N}(0, 1)$
- $v_i \sim \mathcal{U}(\{-1, 1\})$

I guess could be any almost surely **bounded** initialization for \mathbf{v} as we can see later.

Initialization and Conditions

Initialize all the parameters in the network as follows:

- $W_{ij}^{(l)} \overset{i.i.d.}{\sim} \mathcal{N}(0, 1)$
- $v_i \sim \mathcal{U}(\{-1, 1\})$

I guess could be any almost surely **bounded** initialization for \mathbf{v} as we can see later.

Conditions on non-linearity $\sigma : \mathbb{R} \rightarrow \mathbb{R}$:

Initialization and Conditions

Initialize all the parameters in the network as follows:

- $W_{ij}^{(l)} \overset{i.i.d.}{\sim} \mathcal{N}(0, 1)$
- $v_i \sim \mathcal{U}(\{-1, 1\})$

I guess could be any almost surely **bounded** initialization for \mathbf{v} as we can see later.

Conditions on non-linearity $\sigma : \mathbb{R} \rightarrow \mathbb{R}$:

- β -smooth: $\|\sigma(\mathbf{w}) - \sigma(\mathbf{v}) - \nabla\sigma(\mathbf{v})(\mathbf{w} - \mathbf{v})\|_2 \leq \frac{\beta}{2}\|\mathbf{w} - \mathbf{v}\|_2^2$

Initialization and Conditions

Initialize all the parameters in the network as follows:

- $W_{ij}^{(l)} \overset{i.i.d.}{\sim} \mathcal{N}(0, 1)$
- $v_i \sim \mathcal{U}(\{-1, 1\})$

I guess could be any almost surely **bounded** initialization for \mathbf{v} as we can see later.

Conditions on non-linearity $\sigma : \mathbb{R} \rightarrow \mathbb{R}$:

- β -smooth: $\|\sigma(\mathbf{w}) - \sigma(\mathbf{v}) - \nabla\sigma(\mathbf{v})(\mathbf{w} - \mathbf{v})\|_2 \leq \frac{\beta}{2}\|\mathbf{w} - \mathbf{v}\|_2^2$
- Lipschitz: $\|\sigma(\mathbf{w}) - \sigma(\mathbf{v})\|_2 \leq \gamma\|\mathbf{v} - \mathbf{w}\|_2$

Notations

Notations

- We will mostly only deal with the "empirical" NTK (meaning finite width):

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}') = (\nabla_{\mathbf{w}} f(\mathbf{w}; \mathbf{x}))^T \nabla_{\mathbf{w}} f(\mathbf{w}; \mathbf{x}')$$

Notations

- We will mostly only deal with the "empirical" NTK (meaning finite width):

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}') = (\nabla_{\mathbf{w}} f(\mathbf{w}; \mathbf{x}))^T \nabla_{\mathbf{w}} f(\mathbf{w}; \mathbf{x}')$$

- Will view this kernel mostly as a function of \mathbf{w}

Notations

- We will mostly only deal with the "empirical" NTK (meaning finite width):

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}') = (\nabla_{\mathbf{w}} f(\mathbf{w}; \mathbf{x}))^T \nabla_{\mathbf{w}} f(\mathbf{w}; \mathbf{x}')$$

- Will view this kernel mostly as a function of \mathbf{w}
- Of particular interest are diagonal entries:

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \|\nabla_{\mathbf{w}} f(\mathbf{w}; \mathbf{x})\|_2^2$$

Observation I

Observation I

A rather "intuitive" but very important fact is the following:

Observation I

A rather "intuitive" but very important fact is the following:

Constancy \iff Linear

The tangent kernel $\mathbf{w} \mapsto K_{\mathbf{w}}$ of a differentiable model $f(\mathbf{w}; \mathbf{x})$ is **constant** if and only if $f(\mathbf{w}; \mathbf{x})$ is **linear** in \mathbf{w} .

Observation I

A rather "intuitive" but very important fact is the following:

Constancy \iff Linear

The tangent kernel $\mathbf{w} \mapsto K_{\mathbf{w}}$ of a differentiable model $f(\mathbf{w}; \mathbf{x})$ is **constant** if and only if $f(\mathbf{w}; \mathbf{x})$ is **linear** in \mathbf{w} .

Proof is surprisingly long but uses very elementary tools from calculus.

Observation II

Observation II

A new proof strategy for the constancy of NTK:

Observation II

A new proof strategy for the constancy of NTK:

- **Originally:** $\xRightarrow{\text{Jacot}}$ $K_{\mathbf{w}}$ constant $\xRightarrow{\text{Lee et al.}}$ $f(\mathbf{w}; \mathbf{x})$ linear

Observation II

A new proof strategy for the constancy of NTK:

- **Originally:** $\xRightarrow{\text{Jacot}} K_{\mathbf{w}}$ constant $\xRightarrow{\text{Lee et al.}} f(\mathbf{w}; \mathbf{x})$ linear
- **Now:** $\implies f(\mathbf{w}, \mathbf{x})$ linear $\implies K_{\mathbf{w}}$ constant

Observation II

A new proof strategy for the constancy of NTK:

- **Originally:** $\xRightarrow{\text{Jacot}} K_{\mathbf{w}}$ constant $\xRightarrow{\text{Lee et al.}} f(\mathbf{w}; \mathbf{x})$ linear
- **Now:** $\implies f(\mathbf{w}, \mathbf{x})$ linear $\implies K_{\mathbf{w}}$ constant

Let us hence first show that $f(\mathbf{w}; \mathbf{x})$ becomes linear.

Observation III

Observation III

We have the following well-known characterization of linearity:

Observation III

We have the following well-known characterization of linearity:

Linearity vs Hessian

A differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, $\mathbf{w} \mapsto f(\mathbf{w})$ is linear if and only if the Hessian $\mathbf{H}_f = \frac{\partial^2 f}{\partial \mathbf{w}^2}$ vanishes $\forall \mathbf{w} \in \mathbb{R}^d$.

Observation III

We have the following well-known characterization of linearity:

Linearity vs Hessian

A differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, $\mathbf{w} \mapsto f(\mathbf{w})$ is linear if and only if the Hessian $\mathbf{H}_f = \frac{\partial^2 f}{\partial \mathbf{w}^2}$ vanishes $\forall \mathbf{w} \in \mathbb{R}^d$.

It hence suffices to show that $\mathbf{H}_f = 0$ for $m \rightarrow \infty$ or equivalently $\|\mathbf{H}_f\|_2 = 0$

Observation IV

Observation IV

We can convert above asymptotic facts to the **non-asymptotic** (finite-width) setting:

Observation IV

We can convert above asymptotic facts to the **non-asymptotic** (finite-width) setting:

Small Hessian \implies Small change in K_w

Observation IV

We can convert above asymptotic facts to the **non-asymptotic** (finite-width) setting:

Small Hessian \implies Small change in K_w

Fix $\mathbf{w}_0 \in \mathbb{R}^d$ and $R > 0$ and consider $\mathcal{B}(\mathbf{w}_0, R)$. If it holds that $\|\mathbf{H}_f(\mathbf{w})\|_2 < \epsilon \ \forall \mathbf{w} \in \mathcal{B}(\mathbf{w}_0, R)$, then we have that

Observation IV

We can convert above asymptotic facts to the **non-asymptotic** (finite-width) setting:

Small Hessian \implies Small change in K_w

Fix $\mathbf{w}_0 \in \mathbb{R}^d$ and $R > 0$ and consider $\mathcal{B}(\mathbf{w}_0, R)$. If it holds that $\|\mathbf{H}_f(\mathbf{w})\|_2 < \epsilon \ \forall \mathbf{w} \in \mathcal{B}(\mathbf{w}_0, R)$, then we have that

$$|K_w(\mathbf{x}, \mathbf{x}') - K_{\mathbf{w}_0}(\mathbf{x}, \mathbf{x}')| = \mathcal{O}(R\epsilon)$$

$\forall \mathbf{x}, \mathbf{x}' \in \mathbb{R}^{m_0}$ and $\forall \mathbf{w} \in \mathcal{B}(\mathbf{w}_0, R)$

Observation IV

We can convert above asymptotic facts to the **non-asymptotic** (finite-width) setting:

Small Hessian \implies Small change in K_w

Fix $\mathbf{w}_0 \in \mathbb{R}^d$ and $R > 0$ and consider $\mathcal{B}(\mathbf{w}_0, R)$. If it holds that $\|\mathbf{H}_f(\mathbf{w})\|_2 < \epsilon \ \forall \mathbf{w} \in \mathcal{B}(\mathbf{w}_0, R)$, then we have that

$$|K_w(\mathbf{x}, \mathbf{x}') - K_{\mathbf{w}_0}(\mathbf{x}, \mathbf{x}')| = \mathcal{O}(R\epsilon)$$

$$\forall \mathbf{x}, \mathbf{x}' \in \mathbb{R}^{m_0} \text{ and } \forall \mathbf{w} \in \mathcal{B}(\mathbf{w}_0, R)$$

Hence if $\epsilon = \epsilon(m) \xrightarrow{m \rightarrow \infty} 0$ then also

$$|K_w(\mathbf{x}, \mathbf{x}') - K_{\mathbf{w}_0}(\mathbf{x}, \mathbf{x}')| \xrightarrow{m \rightarrow \infty} 0$$

Observation V

Observation V

- **Notice:** We only have a valid result in $\mathcal{B}(\mathbf{w}_0, R)$

Observation V

- **Notice:** We only have a valid result in $\mathcal{B}(\mathbf{w}_0, R)$
- We don't want constancy $\forall \mathbf{w} \in \mathbb{R}^d$, we only want it along the **training trajectory** of GD!

Observation V

- **Notice:** We only have a valid result in $\mathcal{B}(\mathbf{w}_0, R)$
- We don't want constancy $\forall \mathbf{w} \in \mathbb{R}^d$, we only want it along the **training trajectory** of GD!

Exploration of GD

Observation V

- **Notice:** We only have a valid result in $\mathcal{B}(\mathbf{w}_0, R)$
- We don't want constancy $\forall \mathbf{w} \in \mathbb{R}^d$, we only want it along the **training trajectory** of GD!

Exploration of GD

Take NN $f(\mathbf{w}; \mathbf{x})$ and consider β -smooth and $\mu - PL^*$ loss $L(\mathbf{w})$. Fix $R = \frac{2\sqrt{2\beta L(\mathbf{w}_0)}}{\mu}$. Then it holds that:

Observation V

- **Notice:** We only have a valid result in $\mathcal{B}(\mathbf{w}_0, R)$
- We don't want constancy $\forall \mathbf{w} \in \mathbb{R}^d$, we only want it along the **training trajectory** of GD!

Exploration of GD

Take NN $f(\mathbf{w}; \mathbf{x})$ and consider β -smooth and $\mu - PL^*$ loss $L(\mathbf{w})$. Fix $R = \frac{2\sqrt{2\beta L(\mathbf{w}_0)}}{\mu}$. Then it holds that:

- $\exists \mathbf{w}^* \in \mathcal{B}(\mathbf{w}_0, R)$ such that $L(\mathbf{w}^*) = 0$

Observation V

- **Notice:** We only have a valid result in $\mathcal{B}(\mathbf{w}_0, R)$
- We don't want constancy $\forall \mathbf{w} \in \mathbb{R}^d$, we only want it along the **training trajectory** of GD!

Exploration of GD

Take NN $f(\mathbf{w}; \mathbf{x})$ and consider β -smooth and $\mu - PL^*$ loss $L(\mathbf{w})$. Fix $R = \frac{2\sqrt{2\beta L(\mathbf{w}_0)}}{\mu}$. Then it holds that:

- $\exists \mathbf{w}^* \in \mathcal{B}(\mathbf{w}_0, R)$ such that $L(\mathbf{w}^*) = 0$
- GD with $\eta = \frac{1}{\sup_{\mathbf{w} \in \mathcal{B}(\mathbf{w}_0, R)} \lambda_{\max}(\mathbf{H}_L)}$ converges in $\mathcal{B}(\mathbf{w}_0, R)$

Observation V

- **Notice:** We only have a valid result in $\mathcal{B}(\mathbf{w}_0, R)$
- We don't want constancy $\forall \mathbf{w} \in \mathbb{R}^d$, we only want it along the **training trajectory** of GD!

Exploration of GD

Take NN $f(\mathbf{w}; \mathbf{x})$ and consider β -smooth and $\mu - PL^*$ loss $L(\mathbf{w})$. Fix $R = \frac{2\sqrt{2\beta L(\mathbf{w}_0)}}{\mu}$. Then it holds that:

- $\exists \mathbf{w}^* \in \mathcal{B}(\mathbf{w}_0, R)$ such that $L(\mathbf{w}^*) = 0$
- GD with $\eta = \frac{1}{\sup_{\mathbf{w} \in \mathcal{B}(\mathbf{w}_0, R)} \lambda_{\max}(\mathbf{H}_L)}$ converges in $\mathcal{B}(\mathbf{w}_0, R)$

Thus: Fixing $R = \frac{2\sqrt{2\beta L(\mathbf{w}_0)}}{\mu}$ above (independent of m) **suffices!**

Warm-up: 1 hidden layer

Warm-up: 1 hidden layer

Consider 1 hidden layer network with one dimensional input $x \in \mathbb{R}$:

Warm-up: 1 hidden layer

Consider 1 hidden layer network with one dimensional input $x \in \mathbb{R}$:

- $f(\mathbf{w}; x) = \frac{1}{\sqrt{m}} \sum_{i=1}^m v_i \alpha_i(x)$

Warm-up: 1 hidden layer

Consider 1 hidden layer network with one dimensional input $x \in \mathbb{R}$:

- $f(\mathbf{w}; x) = \frac{1}{\sqrt{m}} \sum_{i=1}^m v_i \alpha_i(x)$
- $\alpha_i(x) = \sigma(w_i x)$

where $\mathbf{w} \in \mathbb{R}^m$ and $\mathbf{v} \in \mathbb{R}^m$.

Hessian H_f

Hessian H_f

Let's calculate the **Hessian** of $f(\mathbf{w}; x)$:

Hessian H_f

Let's calculate the **Hessian** of $f(\mathbf{w}; x)$:

$$H_{ij} = \frac{\partial^2 f}{\partial w_i \partial w_j}$$

Hessian H_f

Let's calculate the **Hessian** of $f(\mathbf{w}; x)$:

$$H_{ij} = \frac{\partial^2 f}{\partial w_i \partial w_j} = \frac{\partial^2}{\partial w_i \partial w_j} \frac{1}{\sqrt{m}} \sum_{k=1}^m v_k \sigma(w_k x)$$

Hessian H_f

Let's calculate the **Hessian** of $f(\mathbf{w}; x)$:

$$H_{ij} = \frac{\partial^2 f}{\partial w_i \partial w_j} = \frac{\partial^2}{\partial w_i \partial w_j} \frac{1}{\sqrt{m}} \sum_{k=1}^m v_k \sigma(w_k x) = \frac{\partial}{\partial w_i} \frac{1}{\sqrt{m}} v_j \sigma'(w_j x) x$$

Hessian H_f

Let's calculate the **Hessian** of $f(\mathbf{w}; x)$:

$$\begin{aligned} H_{ij} &= \frac{\partial^2 f}{\partial w_i \partial w_j} = \frac{\partial^2}{\partial w_i \partial w_j} \frac{1}{\sqrt{m}} \sum_{k=1}^m v_k \sigma(w_k x) = \frac{\partial}{\partial w_i} \frac{1}{\sqrt{m}} v_j \sigma'(w_j x) x \\ &= \frac{1}{\sqrt{m}} v_j x^2 \sigma''(w_j x) \mathbb{1}_{\{i=j\}} \end{aligned}$$

Hessian H_f

Let's calculate the **Hessian** of $f(\mathbf{w}; x)$:

$$\begin{aligned} H_{ij} &= \frac{\partial^2 f}{\partial w_i \partial w_j} = \frac{\partial^2}{\partial w_i \partial w_j} \frac{1}{\sqrt{m}} \sum_{k=1}^m v_k \sigma(w_k x) = \frac{\partial}{\partial w_i} \frac{1}{\sqrt{m}} v_j \sigma'(w_j x) x \\ &= \frac{1}{\sqrt{m}} v_j x^2 \sigma''(w_j x) \mathbb{1}_{\{i=j\}} \end{aligned}$$

Since \mathbf{H} is diagonal:

Hessian H_f

Let's calculate the **Hessian** of $f(\mathbf{w}; x)$:

$$\begin{aligned} H_{ij} &= \frac{\partial^2 f}{\partial w_i \partial w_j} = \frac{\partial^2}{\partial w_i \partial w_j} \frac{1}{\sqrt{m}} \sum_{k=1}^m v_k \sigma(w_k x) = \frac{\partial}{\partial w_i} \frac{1}{\sqrt{m}} v_j \sigma'(w_j x) x \\ &= \frac{1}{\sqrt{m}} v_j x^2 \sigma''(w_j x) \mathbb{1}_{\{i=j\}} \end{aligned}$$

Since \mathbf{H} is diagonal:

$$\|\mathbf{H}\|_2 = \sigma_{\max}(\mathbf{H})$$

Hessian H_f

Let's calculate the **Hessian** of $f(\mathbf{w}; x)$:

$$\begin{aligned} H_{ij} &= \frac{\partial^2 f}{\partial w_i \partial w_j} = \frac{\partial^2}{\partial w_i \partial w_j} \frac{1}{\sqrt{m}} \sum_{k=1}^m v_k \sigma(w_k x) = \frac{\partial}{\partial w_i} \frac{1}{\sqrt{m}} v_j \sigma'(w_j x) x \\ &= \frac{1}{\sqrt{m}} v_j x^2 \sigma''(w_j x) \mathbb{1}_{\{i=j\}} \end{aligned}$$

Since \mathbf{H} is diagonal:

$$\|\mathbf{H}\|_2 = \sigma_{\max}(\mathbf{H}) = \max_{j=1}^m |H_{jj}|$$

Hessian H_f

Let's calculate the **Hessian** of $f(\mathbf{w}; x)$:

$$\begin{aligned} H_{ij} &= \frac{\partial^2 f}{\partial w_i \partial w_j} = \frac{\partial^2}{\partial w_i \partial w_j} \frac{1}{\sqrt{m}} \sum_{k=1}^m v_k \sigma(w_k x) = \frac{\partial}{\partial w_i} \frac{1}{\sqrt{m}} v_j \sigma'(w_j x) x \\ &= \frac{1}{\sqrt{m}} v_j x^2 \sigma''(w_j x) \mathbb{1}_{\{i=j\}} \end{aligned}$$

Since \mathbf{H} is diagonal:

$$\|\mathbf{H}\|_2 = \sigma_{\max}(\mathbf{H}) = \max_{j=1}^m |H_{jj}| = \max_{j=1}^m \frac{1}{\sqrt{m}} x^2 |v_j| |\sigma''(w_j x)|$$

Hessian H_f

Let's calculate the **Hessian** of $f(\mathbf{w}; x)$:

$$\begin{aligned} H_{ij} &= \frac{\partial^2 f}{\partial w_i \partial w_j} = \frac{\partial^2}{\partial w_i \partial w_j} \frac{1}{\sqrt{m}} \sum_{k=1}^m v_k \sigma(w_k x) = \frac{\partial}{\partial w_i} \frac{1}{\sqrt{m}} v_j \sigma'(w_j x) x \\ &= \frac{1}{\sqrt{m}} v_j x^2 \sigma''(w_j x) \mathbb{1}_{\{i=j\}} \end{aligned}$$

Since \mathbf{H} is diagonal:

$$\begin{aligned} \|\mathbf{H}\|_2 &= \sigma_{\max}(\mathbf{H}) = \max_{j=1}^m |H_{jj}| = \max_{j=1}^m \frac{1}{\sqrt{m}} x^2 |v_j| |\sigma''(w_j x)| \\ &\leq \frac{1}{\sqrt{m}} C^2 \beta_\sigma \end{aligned}$$

Hessian H_f

Let's calculate the **Hessian** of $f(\mathbf{w}; x)$:

$$\begin{aligned} H_{ij} &= \frac{\partial^2 f}{\partial w_i \partial w_j} = \frac{\partial^2}{\partial w_i \partial w_j} \frac{1}{\sqrt{m}} \sum_{k=1}^m v_k \sigma(w_k x) = \frac{\partial}{\partial w_i} \frac{1}{\sqrt{m}} v_j \sigma'(w_j x) x \\ &= \frac{1}{\sqrt{m}} v_j x^2 \sigma''(w_j x) \mathbb{1}_{\{i=j\}} \end{aligned}$$

Since \mathbf{H} is diagonal:

$$\begin{aligned} \|\mathbf{H}\|_2 &= \sigma_{\max}(\mathbf{H}) = \max_{j=1}^m |H_{jj}| = \max_{j=1}^m \frac{1}{\sqrt{m}} x^2 |v_j| |\sigma''(w_j x)| \\ &\leq \frac{1}{\sqrt{m}} C^2 \beta_\sigma = \mathcal{O}\left(\frac{1}{\sqrt{m}}\right) \end{aligned}$$

Hessian H_f

Let's calculate the **Hessian** of $f(\mathbf{w}; x)$:

$$\begin{aligned} H_{ij} &= \frac{\partial^2 f}{\partial w_i \partial w_j} = \frac{\partial^2}{\partial w_i \partial w_j} \frac{1}{\sqrt{m}} \sum_{k=1}^m v_k \sigma(w_k x) = \frac{\partial}{\partial w_i} \frac{1}{\sqrt{m}} v_j \sigma'(w_j x) x \\ &= \frac{1}{\sqrt{m}} v_j x^2 \sigma''(w_j x) \mathbb{1}_{\{i=j\}} \end{aligned}$$

Since \mathbf{H} is diagonal:

$$\begin{aligned} \|\mathbf{H}\|_2 &= \sigma_{\max}(\mathbf{H}) = \max_{j=1}^m |H_{jj}| = \max_{j=1}^m \frac{1}{\sqrt{m}} x^2 |v_j| |\sigma''(w_j x)| \\ &\leq \frac{1}{\sqrt{m}} C^2 \beta_\sigma = \mathcal{O}\left(\frac{1}{\sqrt{m}}\right) \end{aligned}$$

Thus: $\|\mathbf{H}\|_2 \xrightarrow{m \rightarrow \infty} 0$

Diagonal Kernel Entry $K_w(x, x)$

Diagonal Kernel Entry $K_w(x, x)$

On the other hand we have that:

Diagonal Kernel Entry $K_w(x, x)$

On the other hand we have that:

$$\frac{\partial f}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{\sqrt{m}} \sum_{k=1}^m v_k \sigma(w_k x)$$

Diagonal Kernel Entry $K_w(x, x)$

On the other hand we have that:

$$\frac{\partial f}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{\sqrt{m}} \sum_{k=1}^m v_k \sigma(w_k x) = \frac{1}{\sqrt{m}} v_i \sigma'(w_i x) x$$

Diagonal Kernel Entry $K_w(\mathbf{x}, \mathbf{x})$

On the other hand we have that:

$$\frac{\partial f}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{\sqrt{m}} \sum_{k=1}^m v_k \sigma(w_k x) = \frac{1}{\sqrt{m}} v_i \sigma'(w_i x) x$$

Thus:

$$\|\nabla f(\mathbf{w}, x)\|_2^2 = \frac{x^2}{m} \sum_{i=1}^m v_i^2 \sigma'(w_i x)$$

Diagonal Kernel Entry $K_w(\mathbf{x}, \mathbf{x})$

On the other hand we have that:

$$\frac{\partial f}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{\sqrt{m}} \sum_{k=1}^m v_k \sigma(w_k x) = \frac{1}{\sqrt{m}} v_i \sigma'(w_i x) x$$

Thus:

$$\|\nabla f(\mathbf{w}, x)\|_2^2 = \frac{x^2}{m} \sum_{i=1}^m v_i^2 \sigma'(w_i x) = \frac{x^2}{m} \sum_{i=1}^m \sigma'(w_i x)$$

Diagonal Kernel Entry $K_w(\mathbf{x}, \mathbf{x})$

On the other hand we have that:

$$\frac{\partial f}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{\sqrt{m}} \sum_{k=1}^m v_k \sigma(w_k x) = \frac{1}{\sqrt{m}} v_i \sigma'(w_i x) x$$

Thus:

$$\|\nabla f(\mathbf{w}, x)\|_2^2 = \frac{x^2}{m} \sum_{i=1}^m v_i^2 \sigma'(w_i x) = \frac{x^2}{m} \sum_{i=1}^m \sigma'(w_i x) = \Theta(1)$$

Diagonal Kernel Entry $K_{\mathbf{w}}(\mathbf{x}, \mathbf{x})$

On the other hand we have that:

$$\frac{\partial f}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{\sqrt{m}} \sum_{k=1}^m v_k \sigma(w_k x) = \frac{1}{\sqrt{m}} v_i \sigma'(w_i x) x$$

Thus:

$$\|\nabla f(\mathbf{w}, \mathbf{x})\|_2^2 = \frac{x^2}{m} \sum_{i=1}^m v_i^2 \sigma'(w_i x) = \frac{x^2}{m} \sum_{i=1}^m \sigma'(w_i x) = \Theta(1)$$

Moreover:

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = x^2 \frac{1}{m} \sum_{i=1}^m \sigma'(w_i x)$$

Diagonal Kernel Entry $K_{\mathbf{w}}(\mathbf{x}, \mathbf{x})$

On the other hand we have that:

$$\frac{\partial f}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{\sqrt{m}} \sum_{k=1}^m v_k \sigma(w_k x) = \frac{1}{\sqrt{m}} v_i \sigma'(w_i x) x$$

Thus:

$$\|\nabla f(\mathbf{w}, x)\|_2^2 = \frac{x^2}{m} \sum_{i=1}^m v_i^2 \sigma'(w_i x) = \frac{x^2}{m} \sum_{i=1}^m \sigma'(w_i x) = \Theta(1)$$

Moreover:

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = x^2 \frac{1}{m} \sum_{i=1}^m \sigma'(w_i x) \xrightarrow{m \rightarrow \infty} x^2 \mathbb{E}_{w \sim \mathcal{N}(0,1)} [\sigma'(w x)]$$

Why Different Scalings?

Why Different Scalings?

- **Question:** How come that different scalings in m arise in **Hessian** norm and **gradient** norm??

Why Different Scalings?

- **Question:** How come that different scalings in m arise in **Hessian** norm and **gradient** norm??
- Let's study the order 3 tensor $\mathbf{T} = \frac{\partial \alpha}{\partial \mathbf{w}^2}$

Why Different Scalings?

- **Question:** How come that different scalings in m arise in **Hessian** norm and **gradient** norm??
- Let's study the order 3 tensor $\mathbf{T} = \frac{\partial \alpha}{\partial \mathbf{w}^2}$
- Equip it with the **norm**

$$\|\mathbf{T}\|_{2,2,1} = \sup_{\|\mathbf{a}\|_2=\|\mathbf{b}\|_2=1} \sum_k \left| \sum_i \sum_j T_{ijk} a_i b_j \right|$$

Why Different Scalings?

- **Question:** How come that different scalings in m arise in **Hessian** norm and **gradient** norm??
- Let's study the order 3 tensor $\mathbf{T} = \frac{\partial \alpha}{\partial \mathbf{w}^2}$
- Equip it with the **norm**

$$\|\mathbf{T}\|_{2,2,1} = \sup_{\|\mathbf{a}\|_2=\|\mathbf{b}\|_2=1} \sum_k \left| \sum_i \sum_j T_{ijk} a_i b_j \right|$$

- We have the following easy consequence of Hölder inequality:

Why Different Scalings?

- **Question:** How come that different scalings in m arise in **Hessian** norm and **gradient** norm??
- Let's study the order 3 tensor $\mathbf{T} = \frac{\partial \alpha}{\partial \mathbf{w}^2}$
- Equip it with the **norm**

$$\|\mathbf{T}\|_{2,2,1} = \sup_{\|\mathbf{a}\|_2=\|\mathbf{b}\|_2=1} \sum_k \left| \sum_i \sum_j T_{ijk} a_i b_j \right|$$

- We have the following easy consequence of Hölder inequality:

Given $\mathbf{T} \in \mathbb{R}^{d_1 \times d_2 \times d_3}$ and $\mathbf{v} \in \mathbb{R}^{d_3}$, consider $A_{ij} = \sum_k T_{ijk} v_k$. It holds that

$$\|\mathbf{A}\|_2 \leq \|\mathbf{T}\|_{2,2,1} \|\mathbf{v}\|_\infty$$

Hessian and $\|\cdot\|_\infty$ (I)

Hessian and $\|\cdot\|_\infty$ (I)

- This time we express \mathbf{H} in terms of \mathbf{T} :

Hessian and $\|\cdot\|_\infty$ (I)

- This time we express \mathbf{H} in terms of \mathbf{T} :

$$H_{ij} = \frac{\partial^2 f}{\partial w_i \partial w_j}$$

Hessian and $\|\cdot\|_\infty$ (I)

- This time we express \mathbf{H} in terms of \mathbf{T} :

$$H_{ij} = \frac{\partial^2 f}{\partial w_i \partial w_j} = \sum_{k=1}^m \frac{\partial f}{\partial \alpha_k} \frac{\partial^2 \alpha_k}{\partial w_i \partial w_j}$$

Hessian and $\|\cdot\|_\infty$ (I)

- This time we express \mathbf{H} in terms of \mathbf{T} :

$$H_{ij} = \frac{\partial^2 f}{\partial w_i \partial w_j} = \sum_{k=1}^m \frac{\partial f}{\partial \alpha_k} \frac{\partial^2 \alpha_k}{\partial w_i \partial w_j} = \sum_{k=1}^m \frac{\partial f}{\partial \alpha_k} T_{ijk}$$

Hessian and $\|\cdot\|_\infty$ (I)

- This time we express \mathbf{H} in terms of \mathbf{T} :

$$H_{ij} = \frac{\partial^2 f}{\partial w_i \partial w_j} = \sum_{k=1}^m \frac{\partial f}{\partial \alpha_k} \frac{\partial^2 \alpha_k}{\partial w_i \partial w_j} = \sum_{k=1}^m \frac{\partial f}{\partial \alpha_k} T_{ijk}$$

- From previous statement:

Hessian and $\|\cdot\|_\infty$ (I)

- This time we express \mathbf{H} in terms of \mathbf{T} :

$$H_{ij} = \frac{\partial^2 f}{\partial w_i \partial w_j} = \sum_{k=1}^m \frac{\partial f}{\partial \alpha_k} \frac{\partial^2 \alpha_k}{\partial w_i \partial w_j} = \sum_{k=1}^m \frac{\partial f}{\partial \alpha_k} T_{ijk}$$

- From previous statement:

$$\|\mathbf{H}\|_2 \leq \|\mathbf{T}\|_{2,2,1} \left\| \frac{\partial f}{\partial \boldsymbol{\alpha}} \right\|_\infty$$

Hessian and $\|\cdot\|_\infty$ (I)

- This time we express \mathbf{H} in terms of \mathbf{T} :

$$H_{ij} = \frac{\partial^2 f}{\partial w_i \partial w_j} = \sum_{k=1}^m \frac{\partial f}{\partial \alpha_k} \frac{\partial^2 \alpha_k}{\partial w_i \partial w_j} = \sum_{k=1}^m \frac{\partial f}{\partial \alpha_k} T_{ijk}$$

- From previous statement:

$$\|\mathbf{H}\|_2 \leq \|\mathbf{T}\|_{2,2,1} \left\| \frac{\partial f}{\partial \boldsymbol{\alpha}} \right\|_\infty$$

- Observe that

Hessian and $\|\cdot\|_\infty$ (I)

- This time we express \mathbf{H} in terms of \mathbf{T} :

$$H_{ij} = \frac{\partial^2 f}{\partial w_i \partial w_j} = \sum_{k=1}^m \frac{\partial f}{\partial \alpha_k} \frac{\partial^2 \alpha_k}{\partial w_i \partial w_j} = \sum_{k=1}^m \frac{\partial f}{\partial \alpha_k} T_{ijk}$$

- From previous statement:

$$\|\mathbf{H}\|_2 \leq \|\mathbf{T}\|_{2,2,1} \left\| \frac{\partial f}{\partial \boldsymbol{\alpha}} \right\|_\infty$$

- Observe that

$$T_{ijk} = \frac{\partial^2 \alpha_k}{\partial w_i \partial w_j}$$

Hessian and $\|\cdot\|_\infty$ (I)

- This time we express \mathbf{H} in terms of \mathbf{T} :

$$H_{ij} = \frac{\partial^2 f}{\partial w_i \partial w_j} = \sum_{k=1}^m \frac{\partial f}{\partial \alpha_k} \frac{\partial^2 \alpha_k}{\partial w_i \partial w_j} = \sum_{k=1}^m \frac{\partial f}{\partial \alpha_k} T_{ijk}$$

- From previous statement:

$$\|\mathbf{H}\|_2 \leq \|\mathbf{T}\|_{2,2,1} \left\| \frac{\partial f}{\partial \boldsymbol{\alpha}} \right\|_\infty$$

- Observe that

$$T_{ijk} = \frac{\partial^2 \alpha_k}{\partial w_i \partial w_j} = \frac{\partial^2}{\partial w_i \partial w_j} \sigma(w_k x)$$

Hessian and $\|\cdot\|_\infty$ (I)

- This time we express \mathbf{H} in terms of \mathbf{T} :

$$H_{ij} = \frac{\partial^2 f}{\partial w_i \partial w_j} = \sum_{k=1}^m \frac{\partial f}{\partial \alpha_k} \frac{\partial^2 \alpha_k}{\partial w_i \partial w_j} = \sum_{k=1}^m \frac{\partial f}{\partial \alpha_k} T_{ijk}$$

- From previous statement:

$$\|\mathbf{H}\|_2 \leq \|\mathbf{T}\|_{2,2,1} \left\| \frac{\partial f}{\partial \boldsymbol{\alpha}} \right\|_\infty$$

- Observe that

$$T_{ijk} = \frac{\partial^2 \alpha_k}{\partial w_i \partial w_j} = \frac{\partial^2}{\partial w_i \partial w_j} \sigma(w_k x) = \sigma''(w_k x) x^2 \mathbb{1}_{\{i=j=k\}}$$

Hessian and $\|\cdot\|_\infty$ (II)

Hessian and $\|\cdot\|_\infty$ (II)

- We can hence calculate the norm as:

Hessian and $\|\cdot\|_\infty$ (II)

- We can hence calculate the norm as:

$$\|\mathbf{T}\|_{2,2,1} = \sup_k \left| \sum_i \sum_j T_{ijl} a_i b_j \right|$$

Hessian and $\|\cdot\|_\infty$ (II)

- We can hence calculate the norm as:

$$\|\mathbf{T}\|_{2,2,1} = \sup_k \sum_k \left| \sum_i \sum_j T_{ijl} a_i b_i \right| = x^2 \sup \sum_{k=1}^m |\sigma''(w_k x)|$$

Hessian and $\|\cdot\|_\infty$ (II)

- We can hence calculate the norm as:

$$\begin{aligned}\|\mathbf{T}\|_{2,2,1} &= \sup_k \sum_k \left| \sum_i \sum_j T_{ijl} a_i b_j \right| = x^2 \sup \sum_{k=1}^m |\sigma''(w_k x)| \\ &\leq x^2 L_\sigma \sup_{\|\mathbf{a}\|_2 = \|\mathbf{b}\|_2 = 1} \sum_{k=1} |a_i b_j|\end{aligned}$$

Hessian and $\|\cdot\|_\infty$ (II)

- We can hence calculate the norm as:

$$\begin{aligned}\|\mathbf{T}\|_{2,2,1} &= \sup_k \sum_i \left| \sum_j T_{ijl} a_i b_j \right| = x^2 \sup_k \sum_{i=1}^m |\sigma''(w_k x)| \\ &\leq x^2 L_\sigma \sup_{\|\mathbf{a}\|_2=\|\mathbf{b}\|_2=1} \sum_{i=1}^m |a_i b_i| \leq x^2 L_\sigma\end{aligned}$$

Hessian and $\|\cdot\|_\infty$ (II)

- We can hence calculate the norm as:

$$\begin{aligned}\|\mathbf{T}\|_{2,2,1} &= \sup_k \sum_i \left| \sum_j T_{ijl} a_i b_j \right| = x^2 \sup_k \sum_{i=1}^m |\sigma''(w_k x)| \\ &\leq x^2 L_\sigma \sup_{\|\mathbf{a}\|_2=\|\mathbf{b}\|_2=1} \sum_{i=1}^m |a_i b_i| \leq x^2 L_\sigma = \mathcal{O}(1)\end{aligned}$$

Hessian and $\|\cdot\|_\infty$ (II)

- We can hence calculate the norm as:

$$\begin{aligned}\|\mathbf{T}\|_{2,2,1} &= \sup_k \sum_i \left| \sum_j T_{ijl} a_i b_j \right| = x^2 \sup_k \sum_{i=1}^m |\sigma''(w_k x)| \\ &\leq x^2 L_\sigma \sup_{\|\mathbf{a}\|_2=\|\mathbf{b}\|_2=1} \sum_{i=1}^m |a_i b_i| \leq x^2 L_\sigma = \mathcal{O}(1)\end{aligned}$$

- This in turn implies that

Hessian and $\|\cdot\|_\infty$ (II)

- We can hence calculate the norm as:

$$\begin{aligned}\|\mathbf{T}\|_{2,2,1} &= \sup_k \sum_i \left| \sum_j T_{ijl} a_i b_j \right| = x^2 \sup_k \sum_{i=1}^m |\sigma''(w_k x)| \\ &\leq x^2 L_\sigma \sup_{\|\mathbf{a}\|_2=\|\mathbf{b}\|_2=1} \sum_{i=1}^m |a_i b_i| \leq x^2 L_\sigma = \mathcal{O}(1)\end{aligned}$$

- This in turn implies that

$$\|\mathbf{H}\|_2 = \mathcal{O}\left(\left\|\frac{\partial f}{\partial \boldsymbol{\alpha}}\right\|_\infty\right)$$

Hessian and $\|\cdot\|_\infty$ (II)

- We can hence calculate the norm as:

$$\begin{aligned}\|\mathbf{T}\|_{2,2,1} &= \sup_k \sum_i \left| \sum_j T_{ijl} a_i b_j \right| = x^2 \sup_{k=1}^m |\sigma''(w_k x)| \\ &\leq x^2 L_\sigma \sup_{\|\mathbf{a}\|_2=\|\mathbf{b}\|_2=1} \sum_{k=1} |a_i b_j| \leq x^2 L_\sigma = \mathcal{O}(1)\end{aligned}$$

- This in turn implies that

$$\|\mathbf{H}\|_2 = \mathcal{O}\left(\left\|\frac{\partial f}{\partial \boldsymbol{\alpha}}\right\|_\infty\right)$$

- The l_∞ - norm thus governs the scaling of the Hessian

Kernel and $\|\cdot\|_2$

Kernel and $\|\cdot\|_2$

- Observe the following for the kernel $\mathbf{K} \in \mathbb{R}^{n \times n}$:

Kernel and $\|\cdot\|_2$

- Observe the following for the kernel $\mathbf{K} \in \mathbb{R}^{n \times n}$:

$$\|\mathbf{K}\|_2 = \lambda_{\max}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{K})$$

Kernel and $\|\cdot\|_2$

- Observe the following for the kernel $\mathbf{K} \in \mathbb{R}^{n \times n}$:

$$\|\mathbf{K}\|_2 = \lambda_{\max}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{K}) \geq \frac{1}{n} \sum_{i=1}^n \lambda_i(\mathbf{K})$$

Kernel and $\|\cdot\|_2$

- Observe the following for the kernel $\mathbf{K} \in \mathbb{R}^{n \times n}$:

$$\begin{aligned}\|\mathbf{K}\|_2 &= \lambda_{\max}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{K}) \geq \frac{1}{n} \sum_{i=1}^n \lambda_i(\mathbf{K}) \\ &= \frac{1}{n} \text{tr}(\mathbf{K})\end{aligned}$$

Kernel and $\|\cdot\|_2$

- Observe the following for the kernel $\mathbf{K} \in \mathbb{R}^{n \times n}$:

$$\begin{aligned}\|\mathbf{K}\|_2 &= \lambda_{\max}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{K}) \geq \frac{1}{n} \sum_{i=1}^n \lambda_i(\mathbf{K}) \\ &= \frac{1}{n} \text{tr}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n K_{\mathbf{w}}(\mathbf{x}_i, \mathbf{x}_i)\end{aligned}$$

Kernel and $\|\cdot\|_2$

- Observe the following for the kernel $\mathbf{K} \in \mathbb{R}^{n \times n}$:

$$\begin{aligned}\|\mathbf{K}\|_2 &= \lambda_{\max}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{K}) \geq \frac{1}{n} \sum_{i=1}^n \lambda_i(\mathbf{K}) \\ &= \frac{1}{n} \text{tr}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n K_{\mathbf{w}}(\mathbf{x}_i, \mathbf{x}_i)\end{aligned}$$

- Consider thus a diagonal kernel entry:

Kernel and $\|\cdot\|_2$

- Observe the following for the kernel $\mathbf{K} \in \mathbb{R}^{n \times n}$:

$$\begin{aligned}\|\mathbf{K}\|_2 &= \lambda_{\max}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{K}) \geq \frac{1}{n} \sum_{i=1}^n \lambda_i(\mathbf{K}) \\ &= \frac{1}{n} \text{tr}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n K_{\mathbf{w}}(\mathbf{x}_i, \mathbf{x}_i)\end{aligned}$$

- Consider thus a diagonal kernel entry:

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \|\nabla_{\mathbf{w}} f(\mathbf{w}; \mathbf{x})\|_2^2$$

Kernel and $\|\cdot\|_2$

- Observe the following for the kernel $\mathbf{K} \in \mathbb{R}^{n \times n}$:

$$\begin{aligned}\|\mathbf{K}\|_2 &= \lambda_{\max}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{K}) \geq \frac{1}{n} \sum_{i=1}^n \lambda_i(\mathbf{K}) \\ &= \frac{1}{n} \text{tr}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n K_{\mathbf{w}}(\mathbf{x}_i, \mathbf{x}_i)\end{aligned}$$

- Consider thus a diagonal kernel entry:

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \|\nabla_{\mathbf{w}} f(\mathbf{w}; \mathbf{x})\|_2^2 = \left\| \frac{\partial \alpha}{\partial \mathbf{w}} \frac{\partial f}{\partial \alpha} \right\|_2^2$$

Kernel and $\|\cdot\|_2$

- Observe the following for the kernel $\mathbf{K} \in \mathbb{R}^{n \times n}$:

$$\begin{aligned}\|\mathbf{K}\|_2 &= \lambda_{\max}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{K}) \geq \frac{1}{n} \sum_{i=1}^n \lambda_i(\mathbf{K}) \\ &= \frac{1}{n} \text{tr}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n K_{\mathbf{w}}(\mathbf{x}_i, \mathbf{x}_i)\end{aligned}$$

- Consider thus a diagonal kernel entry:

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \|\nabla_{\mathbf{w}} f(\mathbf{w}; \mathbf{x})\|_2^2 = \left\| \frac{\partial \alpha}{\partial \mathbf{w}} \frac{\partial f}{\partial \alpha} \right\|_2^2 \approx \left\| \frac{\partial \alpha}{\partial \mathbf{w}} \right\|_2^2 \left\| \frac{\partial f}{\partial \alpha} \right\|_2^2$$

Kernel and $\|\cdot\|_2$

- Observe the following for the kernel $\mathbf{K} \in \mathbb{R}^{n \times n}$:

$$\begin{aligned}\|\mathbf{K}\|_2 &= \lambda_{\max}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{K}) \geq \frac{1}{n} \sum_{i=1}^n \lambda_i(\mathbf{K}) \\ &= \frac{1}{n} \text{tr}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n K_{\mathbf{w}}(\mathbf{x}_i, \mathbf{x}_i)\end{aligned}$$

- Consider thus a diagonal kernel entry:

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \|\nabla_{\mathbf{w}} f(\mathbf{w}; \mathbf{x})\|_2^2 = \left\| \frac{\partial \alpha}{\partial \mathbf{w}} \frac{\partial f}{\partial \alpha} \right\|_2^2 \approx \left\| \frac{\partial \alpha}{\partial \mathbf{w}} \right\|_2^2 \left\| \frac{\partial f}{\partial \alpha} \right\|_2^2$$

- Moreover $\frac{\partial \alpha_i}{\partial w_j}$

Kernel and $\|\cdot\|_2$

- Observe the following for the kernel $\mathbf{K} \in \mathbb{R}^{n \times n}$:

$$\begin{aligned}\|\mathbf{K}\|_2 &= \lambda_{\max}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{K}) \geq \frac{1}{n} \sum_{i=1}^n \lambda_i(\mathbf{K}) \\ &= \frac{1}{n} \text{tr}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n K_{\mathbf{w}}(\mathbf{x}_i, \mathbf{x}_i)\end{aligned}$$

- Consider thus a diagonal kernel entry:

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \|\nabla_{\mathbf{w}} f(\mathbf{w}; \mathbf{x})\|_2^2 = \left\| \frac{\partial \alpha}{\partial \mathbf{w}} \frac{\partial f}{\partial \alpha} \right\|_2^2 \approx \left\| \frac{\partial \alpha}{\partial \mathbf{w}} \right\|_2^2 \left\| \frac{\partial f}{\partial \alpha} \right\|_2^2$$

- Moreover** $\frac{\partial \alpha_i}{\partial w_j} = \sigma'(w_i x) x \mathbb{1}_{\{i=j\}}$

Kernel and $\|\cdot\|_2$

- Observe the following for the kernel $\mathbf{K} \in \mathbb{R}^{n \times n}$:

$$\begin{aligned}\|\mathbf{K}\|_2 &= \lambda_{\max}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{K}) \geq \frac{1}{n} \sum_{i=1}^n \lambda_i(\mathbf{K}) \\ &= \frac{1}{n} \text{tr}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n K_{\mathbf{w}}(\mathbf{x}_i, \mathbf{x}_i)\end{aligned}$$

- Consider thus a diagonal kernel entry:

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \|\nabla_{\mathbf{w}} f(\mathbf{w}; \mathbf{x})\|_2^2 = \left\| \frac{\partial \alpha}{\partial \mathbf{w}} \frac{\partial f}{\partial \alpha} \right\|_2^2 \approx \left\| \frac{\partial \alpha}{\partial \mathbf{w}} \right\|_2^2 \left\| \frac{\partial f}{\partial \alpha} \right\|_2^2$$

- Moreover** $\frac{\partial \alpha_i}{\partial w_j} = \sigma'(w_i x) x \mathbb{1}_{\{i=j\}} \leq L_{\sigma} x$

Kernel and $\|\cdot\|_2$

- Observe the following for the kernel $\mathbf{K} \in \mathbb{R}^{n \times n}$:

$$\begin{aligned}\|\mathbf{K}\|_2 &= \lambda_{\max}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{K}) \geq \frac{1}{n} \sum_{i=1}^n \lambda_i(\mathbf{K}) \\ &= \frac{1}{n} \text{tr}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n K_{\mathbf{w}}(\mathbf{x}_i, \mathbf{x}_i)\end{aligned}$$

- Consider thus a diagonal kernel entry:

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \|\nabla_{\mathbf{w}} f(\mathbf{w}; \mathbf{x})\|_2^2 = \left\| \frac{\partial \boldsymbol{\alpha}}{\partial \mathbf{w}} \frac{\partial f}{\partial \boldsymbol{\alpha}} \right\|_2^2 \approx \left\| \frac{\partial \boldsymbol{\alpha}}{\partial \mathbf{w}} \right\|_2^2 \left\| \frac{\partial f}{\partial \boldsymbol{\alpha}} \right\|_2^2$$

- Moreover** $\frac{\partial \alpha_i}{\partial w_j} = \sigma'(w_i x) x \mathbb{1}_{\{i=j\}} \leq L_{\sigma} x \implies \left\| \frac{\partial \boldsymbol{\alpha}}{\partial \mathbf{w}} \right\|_2^2$

Kernel and $\|\cdot\|_2$

- Observe the following for the kernel $\mathbf{K} \in \mathbb{R}^{n \times n}$:

$$\begin{aligned}\|\mathbf{K}\|_2 &= \lambda_{\max}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{K}) \geq \frac{1}{n} \sum_{i=1}^n \lambda_i(\mathbf{K}) \\ &= \frac{1}{n} \text{tr}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n K_{\mathbf{w}}(\mathbf{x}_i, \mathbf{x}_i)\end{aligned}$$

- Consider thus a diagonal kernel entry:

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \|\nabla_{\mathbf{w}} f(\mathbf{w}; \mathbf{x})\|_2^2 = \left\| \frac{\partial \alpha}{\partial \mathbf{w}} \frac{\partial f}{\partial \alpha} \right\|_2^2 \approx \left\| \frac{\partial \alpha}{\partial \mathbf{w}} \right\|_2^2 \left\| \frac{\partial f}{\partial \alpha} \right\|_2^2$$

- Moreover** $\frac{\partial \alpha_i}{\partial w_j} = \sigma'(w_i x) x \mathbb{1}_{\{i=j\}} \leq L_{\sigma} x \implies \left\| \frac{\partial \alpha}{\partial \mathbf{w}} \right\|_2^2 = \mathcal{O}(1)$

Kernel and $\|\cdot\|_2$

- Observe the following for the kernel $\mathbf{K} \in \mathbb{R}^{n \times n}$:

$$\begin{aligned}\|\mathbf{K}\|_2 &= \lambda_{\max}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{K}) \geq \frac{1}{n} \sum_{i=1}^n \lambda_i(\mathbf{K}) \\ &= \frac{1}{n} \text{tr}(\mathbf{K}) = \frac{1}{n} \sum_{i=1}^n K_{\mathbf{w}}(\mathbf{x}_i, \mathbf{x}_i)\end{aligned}$$

- Consider thus a diagonal kernel entry:

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \|\nabla_{\mathbf{w}} f(\mathbf{w}; \mathbf{x})\|_2^2 = \left\| \frac{\partial \alpha}{\partial \mathbf{w}} \frac{\partial f}{\partial \alpha} \right\|_2^2 \approx \left\| \frac{\partial \alpha}{\partial \mathbf{w}} \right\|_2^2 \left\| \frac{\partial f}{\partial \alpha} \right\|_2^2$$

- Moreover** $\frac{\partial \alpha_i}{\partial w_j} = \sigma'(w_i x) x \mathbb{1}_{\{i=j\}} \leq L_{\sigma} x \implies \left\| \frac{\partial \alpha}{\partial \mathbf{w}} \right\|_2^2 = \mathcal{O}(1)$
- Hence:** $\|\mathbf{K}\|_2 = \mathcal{O} \left(\left\| \frac{\partial f}{\partial \alpha} \right\|_2^2 \right)$

$\|\cdot\|_\infty$ **versus** $\|\cdot\|_2$

$\|\cdot\|_\infty$ **versus** $\|\cdot\|_2$

- Discrepancy hence stems from the different norms:

$\|\cdot\|_\infty$ **versus** $\|\cdot\|_2$

- Discrepancy hence stems from the different norms:

$$1) \quad \|\mathbf{H}\|_2 \leq \mathcal{O}\left(\left\|\frac{\partial f}{\partial \boldsymbol{\alpha}}\right\|_\infty\right)$$

$\|\cdot\|_\infty$ **versus** $\|\cdot\|_2$

- Discrepancy hence stems from the different norms:

$$1) \quad \|\mathbf{H}\|_2 \leq \mathcal{O}\left(\left\|\frac{\partial f}{\partial \boldsymbol{\alpha}}\right\|_\infty\right)$$

$$2) \quad \|\mathbf{K}\|_2 = \mathcal{O}\left(\left\|\frac{\partial f}{\partial \boldsymbol{\alpha}}\right\|_2^2\right)$$

$\|\cdot\|_\infty$ **versus** $\|\cdot\|_2$

- Discrepancy hence stems from the different norms:

$$1) \|\mathbf{H}\|_2 \leq \mathcal{O}\left(\left\|\frac{\partial f}{\partial \boldsymbol{\alpha}}\right\|_\infty\right)$$

$$2) \|\mathbf{K}\|_2 = \mathcal{O}\left(\left\|\frac{\partial f}{\partial \boldsymbol{\alpha}}\right\|_2^2\right)$$

- Now let's calculate this vector:

$\|\cdot\|_\infty$ **versus** $\|\cdot\|_2$

- Discrepancy hence stems from the different norms:

$$1) \|\mathbf{H}\|_2 \leq \mathcal{O}\left(\left\|\frac{\partial f}{\partial \boldsymbol{\alpha}}\right\|_\infty\right)$$

$$2) \|\mathbf{K}\|_2 = \mathcal{O}\left(\left\|\frac{\partial f}{\partial \boldsymbol{\alpha}}\right\|_2^2\right)$$

- Now let's calculate this vector:

$$\frac{\partial f}{\partial \alpha_i} = \frac{\partial}{\partial \alpha_i} \left(\frac{1}{\sqrt{m}} \sum_{j=1}^m v_j \alpha_j \right)$$

$\|\cdot\|_\infty$ **versus** $\|\cdot\|_2$

- Discrepancy hence stems from the different norms:

$$1) \quad \|\mathbf{H}\|_2 \leq \mathcal{O}\left(\left\|\frac{\partial f}{\partial \boldsymbol{\alpha}}\right\|_\infty\right)$$

$$2) \quad \|\mathbf{K}\|_2 = \mathcal{O}\left(\left\|\frac{\partial f}{\partial \boldsymbol{\alpha}}\right\|_2^2\right)$$

- Now let's calculate this vector:

$$\frac{\partial f}{\partial \alpha_i} = \frac{\partial}{\partial \alpha_i} \left(\frac{1}{\sqrt{m}} \sum_{j=1}^m v_j \alpha_j \right) = \frac{1}{\sqrt{m}} v_i$$

$\|\cdot\|_\infty$ **versus** $\|\cdot\|_2$

- Discrepancy hence stems from the different norms:

$$1) \|\mathbf{H}\|_2 \leq \mathcal{O}\left(\left\|\frac{\partial f}{\partial \boldsymbol{\alpha}}\right\|_\infty\right)$$

$$2) \|\mathbf{K}\|_2 = \mathcal{O}\left(\left\|\frac{\partial f}{\partial \boldsymbol{\alpha}}\right\|_2^2\right)$$

- Now let's calculate this vector:

$$\frac{\partial f}{\partial \alpha_i} = \frac{\partial}{\partial \alpha_i} \left(\frac{1}{\sqrt{m}} \sum_{j=1}^m v_j \alpha_j \right) = \frac{1}{\sqrt{m}} v_i$$

- **Thus:**

$\|\cdot\|_\infty$ **versus** $\|\cdot\|_2$

- Discrepancy hence stems from the different norms:

$$1) \|\mathbf{H}\|_2 \leq \mathcal{O}\left(\left\|\frac{\partial f}{\partial \boldsymbol{\alpha}}\right\|_\infty\right)$$

$$2) \|\mathbf{K}\|_2 = \mathcal{O}\left(\left\|\frac{\partial f}{\partial \boldsymbol{\alpha}}\right\|_2^2\right)$$

- Now let's calculate this vector:

$$\frac{\partial f}{\partial \alpha_i} = \frac{\partial}{\partial \alpha_i} \left(\frac{1}{\sqrt{m}} \sum_{j=1}^m v_j \alpha_j \right) = \frac{1}{\sqrt{m}} v_i$$

- **Thus:**

$$1) \left\|\frac{\partial f}{\partial \boldsymbol{\alpha}}\right\|_\infty = \frac{1}{\sqrt{m}}$$

$\|\cdot\|_\infty$ **versus** $\|\cdot\|_2$

- Discrepancy hence stems from the different norms:

$$1) \|\mathbf{H}\|_2 \leq \mathcal{O}\left(\left\|\frac{\partial f}{\partial \boldsymbol{\alpha}}\right\|_\infty\right)$$

$$2) \|\mathbf{K}\|_2 = \mathcal{O}\left(\left\|\frac{\partial f}{\partial \boldsymbol{\alpha}}\right\|_2^2\right)$$

- Now let's calculate this vector:

$$\frac{\partial f}{\partial \alpha_i} = \frac{\partial}{\partial \alpha_i} \left(\frac{1}{\sqrt{m}} \sum_{j=1}^m v_j \alpha_j \right) = \frac{1}{\sqrt{m}} v_i$$

- **Thus:**

$$1) \left\|\frac{\partial f}{\partial \boldsymbol{\alpha}}\right\|_\infty = \frac{1}{\sqrt{m}}$$

$$2) \left\|\frac{\partial f}{\partial \boldsymbol{\alpha}}\right\|_2^2 = \frac{1}{m} \sum_{i=1}^m v_i^2 = 1$$

Setup, Restated

Setup, Restated

- As always, fully-connected network of L layers mapping to \mathbb{R} :

Setup, Restated

- As always, fully-connected network of L layers mapping to \mathbb{R} :
 - $\alpha^{(0)}(\mathbf{x}) = \mathbf{x}$

Setup, Restated

- As always, fully-connected network of L layers mapping to \mathbb{R} :
 - $\alpha^{(0)}(\mathbf{x}) = \mathbf{x}$
 - $\alpha^{(l)}(\mathbf{x}) = \sigma(\tilde{\alpha}^{(l-1)}(\mathbf{x}))$

Setup, Restated

- As always, fully-connected network of L layers mapping to \mathbb{R} :
 - $\alpha^{(0)}(\mathbf{x}) = \mathbf{x}$
 - $\alpha^{(l)}(\mathbf{x}) = \sigma(\tilde{\alpha}^{(l-1)}(\mathbf{x}))$
 - $\tilde{\alpha}^{(l)}(\mathbf{x}) = \frac{1}{\sqrt{m_{l-1}}} \mathbf{W}^{(l)} \alpha^{(l-1)}(\mathbf{x})$

Setup, Restated

- As always, fully-connected network of L layers mapping to \mathbb{R} :
 - $\alpha^{(0)}(\mathbf{x}) = \mathbf{x}$
 - $\alpha^{(l)}(\mathbf{x}) = \sigma(\tilde{\alpha}^{(l-1)}(\mathbf{x}))$
 - $\tilde{\alpha}^{(l)}(\mathbf{x}) = \frac{1}{\sqrt{m_{l-1}}} \mathbf{W}^{(l)} \alpha^{(l-1)}(\mathbf{x})$

where $\mathbf{W}^{(l)} \in \mathbb{R}^{m_{l-1} \times m_l}$ and $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ some **non-linearity**

Setup, Restated

- As always, fully-connected network of L layers mapping to \mathbb{R} :
 - $\alpha^{(0)}(\mathbf{x}) = \mathbf{x}$
 - $\alpha^{(l)}(\mathbf{x}) = \sigma(\tilde{\alpha}^{(l-1)}(\mathbf{x}))$
 - $\tilde{\alpha}^{(l)}(\mathbf{x}) = \frac{1}{\sqrt{m_{l-1}}} \mathbf{W}^{(l)} \alpha^{(l-1)}(\mathbf{x})$

where $\mathbf{W}^{(l)} \in \mathbb{R}^{m_{l-1} \times m_l}$ and $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ some **non-linearity**

- The final output is given by

$$f(\mathbf{x}) = \frac{1}{\sqrt{m}} \mathbf{v}^T \alpha^{(L)}(\mathbf{x})$$

where $\mathbf{v} \in \mathbb{R}^m$ is **not** trainable

Hessian for General Neural Network

Hessian for General Neural Network

- Turns out that the same behaviour is observed!

Hessian for General Neural Network

- Turns out that the same behaviour is observed!
- Define the following quantities:

Hessian for General Neural Network

- Turns out that the same behaviour is observed!
- Define the following quantities:

$$1) \mathcal{Q}_\infty = \max_{1 \leq i \leq L} \left\| \frac{\partial f}{\partial \alpha^{(i)}} \right\|_\infty$$

Hessian for General Neural Network

- Turns out that the same behaviour is observed!
- Define the following quantities:

$$1) \mathcal{Q}_\infty = \max_{1 \leq i \leq L} \left\| \frac{\partial f}{\partial \alpha^{(i)}} \right\|_\infty$$

$$2) \mathcal{Q}_2 = \max_{1 \leq i \leq L} \left\| \frac{\partial f}{\partial \alpha^{(i)}} \right\|_2$$

Hessian for General Neural Network

- Turns out that the same behaviour is observed!
- Define the following quantities:

$$1) \mathcal{Q}_\infty = \max_{1 \leq i \leq L} \left\| \frac{\partial f}{\partial \boldsymbol{\alpha}^{(i)}} \right\|_\infty$$

$$2) \mathcal{Q}_2 = \max_{1 \leq i \leq L} \left\| \frac{\partial f}{\partial \boldsymbol{\alpha}^{(i)}} \right\|_2$$

$$3) \mathcal{Q}_{2,2,1} = \max_{1 \leq l_1 < l_2 < l_3 \leq L} \left\| \frac{\partial \boldsymbol{\alpha}^{(l_1)}}{\partial \mathbf{w}^{(l_1)}} \right\|_2 \left\| \frac{\partial \boldsymbol{\alpha}^{(l_2)}}{\partial \mathbf{w}^{(l_2)}} \right\|_2 \left\| \frac{\partial^2 \boldsymbol{\alpha}^{(l_3)}}{(\partial \boldsymbol{\alpha}^{(l_3-1)})^2} \right\|_{2,2,1}$$

Hessian for General Neural Network

- Turns out that the same behaviour is observed!
- Define the following quantities:

$$1) \mathcal{Q}_\infty = \max_{1 \leq i \leq L} \left\| \frac{\partial f}{\partial \alpha^{(i)}} \right\|_\infty$$

$$2) \mathcal{Q}_2 = \max_{1 \leq i \leq L} \left\| \frac{\partial f}{\partial \alpha^{(i)}} \right\|_2$$

$$3) \mathcal{Q}_{2,2,1} = \max_{1 \leq l_1 < l_2 < l_3 \leq L} \left\| \frac{\partial \alpha^{(l_1)}}{\partial \mathbf{w}^{(l_1)}} \right\|_2 \left\| \frac{\partial \alpha^{(l_2)}}{\partial \mathbf{w}^{(l_2)}} \right\|_2 \left\| \frac{\partial^2 \alpha^{(l_3)}}{(\partial \alpha^{(l_3-1)})^2} \right\|_{2,2,1}$$

Bound on Hessian

We have the following bound on the norm of the Hessian:

$$\| \mathbf{H} \|_2 \leq C_1 \mathcal{Q}_{2,2,1} \mathcal{Q}_\infty + \frac{1}{\sqrt{m}} C_2 \mathcal{Q}_2$$

Kernel for General Neural Network

Kernel for General Neural Network

- We use the same argument:

Kernel for General Neural Network

- We use the same argument:

$$K_w(\mathbf{x}, \mathbf{x})$$

Kernel for General Neural Network

- We use the same argument:

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \sum_{l=1}^L \|\nabla_{\mathbf{w}^{(l)}} f(\mathbf{w}; \mathbf{x})\|_2^2$$

Kernel for General Neural Network

- We use the same argument:

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \sum_{l=1}^L \|\nabla_{\mathbf{w}^{(l)}} f(\mathbf{w}; \mathbf{x})\|_2^2 = \sum_{l=1}^L \left\| \frac{\partial \alpha^{(l)}}{\partial \mathbf{w}^{(l)}} \frac{\partial f}{\partial \alpha^{(l)}} \right\|_2^2$$

Kernel for General Neural Network

- We use the same argument:

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \sum_{l=1}^L \|\nabla_{\mathbf{w}^{(l)}} f(\mathbf{w}; \mathbf{x})\|_2^2 = \sum_{l=1}^L \left\| \frac{\partial \alpha^{(l)}}{\partial \mathbf{w}^{(l)}} \frac{\partial f}{\partial \alpha^{(l)}} \right\|_2^2$$

- Again since all is random at initialization (a bit shady argument):

Kernel for General Neural Network

- We use the same argument:

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \sum_{l=1}^L \|\nabla_{\mathbf{w}^{(l)}} f(\mathbf{w}; \mathbf{x})\|_2^2 = \sum_{l=1}^L \left\| \frac{\partial \alpha^{(l)}}{\partial \mathbf{w}^{(l)}} \frac{\partial f}{\partial \alpha^{(l)}} \right\|_2^2$$

- Again since all is random at initialization (a bit shady argument):

$$\left\| \frac{\partial \alpha^{(l)}}{\partial \mathbf{w}^{(l)}} \frac{\partial f}{\partial \alpha^{(l)}} \right\|_2^2 \approx \left\| \frac{\partial \alpha^{(l)}}{\partial \mathbf{w}^{(l)}} \right\|_2^2 \left\| \frac{\partial f}{\partial \alpha^{(l)}} \right\|_2^2$$

Kernel for General Neural Network

- We use the same argument:

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \sum_{l=1}^L \|\nabla_{\mathbf{w}^{(l)}} f(\mathbf{w}; \mathbf{x})\|_2^2 = \sum_{l=1}^L \left\| \frac{\partial \alpha^{(l)}}{\partial \mathbf{w}^{(l)}} \frac{\partial f}{\partial \alpha^{(l)}} \right\|_2^2$$

- Again since all is random at initialization (a bit shady argument):

$$\left\| \frac{\partial \alpha^{(l)}}{\partial \mathbf{w}^{(l)}} \frac{\partial f}{\partial \alpha^{(l)}} \right\|_2^2 \approx \left\| \frac{\partial \alpha^{(l)}}{\partial \mathbf{w}^{(l)}} \right\|_2^2 \left\| \frac{\partial f}{\partial \alpha^{(l)}} \right\|_2^2$$

- Moreover we get the same asymptotics:

Kernel for General Neural Network

- We use the same argument:

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \sum_{l=1}^L \|\nabla_{\mathbf{w}^{(l)}} f(\mathbf{w}; \mathbf{x})\|_2^2 = \sum_{l=1}^L \left\| \frac{\partial \alpha^{(l)}}{\partial \mathbf{w}^{(l)}} \frac{\partial f}{\partial \alpha^{(l)}} \right\|_2^2$$

- Again since all is random at initialization (a bit shady argument):

$$\left\| \frac{\partial \alpha^{(l)}}{\partial \mathbf{w}^{(l)}} \frac{\partial f}{\partial \alpha^{(l)}} \right\|_2^2 \approx \left\| \frac{\partial \alpha^{(l)}}{\partial \mathbf{w}^{(l)}} \right\|_2^2 \left\| \frac{\partial f}{\partial \alpha^{(l)}} \right\|_2^2$$

- Moreover we get the same asymptotics:

$$\left\| \frac{\partial \alpha^{(l)}}{\partial \mathbf{w}^{(l)}} \right\|_2^2 = \Theta(1)$$

Kernel for General Neural Network

- We use the same argument:

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \sum_{l=1}^L \|\nabla_{\mathbf{w}^{(l)}} f(\mathbf{w}; \mathbf{x})\|_2^2 = \sum_{l=1}^L \left\| \frac{\partial \alpha^{(l)}}{\partial \mathbf{w}^{(l)}} \frac{\partial f}{\partial \alpha^{(l)}} \right\|_2^2$$

- Again since all is random at initialization (a bit shady argument):

$$\left\| \frac{\partial \alpha^{(l)}}{\partial \mathbf{w}^{(l)}} \frac{\partial f}{\partial \alpha^{(l)}} \right\|_2^2 \approx \left\| \frac{\partial \alpha^{(l)}}{\partial \mathbf{w}^{(l)}} \right\|_2^2 \left\| \frac{\partial f}{\partial \alpha^{(l)}} \right\|_2^2$$

- Moreover we get the same asymptotics:

$$\left\| \frac{\partial \alpha^{(l)}}{\partial \mathbf{w}^{(l)}} \right\|_2^2 = \Theta(1)$$

- **So:** $K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \mathcal{O} \left(\max_{l=1}^L \left\| \frac{\partial f}{\partial \alpha^{(l)}} \right\|_2^2 \right)$

Kernel for General Neural Network

- We use the same argument:

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \sum_{l=1}^L \|\nabla_{\mathbf{w}^{(l)}} f(\mathbf{w}; \mathbf{x})\|_2^2 = \sum_{l=1}^L \left\| \frac{\partial \alpha^{(l)}}{\partial \mathbf{w}^{(l)}} \frac{\partial f}{\partial \alpha^{(l)}} \right\|_2^2$$

- Again since all is random at initialization (a bit shady argument):

$$\left\| \frac{\partial \alpha^{(l)}}{\partial \mathbf{w}^{(l)}} \frac{\partial f}{\partial \alpha^{(l)}} \right\|_2^2 \approx \left\| \frac{\partial \alpha^{(l)}}{\partial \mathbf{w}^{(l)}} \right\|_2^2 \left\| \frac{\partial f}{\partial \alpha^{(l)}} \right\|_2^2$$

- Moreover we get the same asymptotics:

$$\left\| \frac{\partial \alpha^{(l)}}{\partial \mathbf{w}^{(l)}} \right\|_2^2 = \Theta(1)$$

- **So:** $K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \mathcal{O} \left(\max_{i=1}^L \left\| \frac{\partial f}{\partial \alpha^{(i)}} \right\|_2^2 \right) = \mathcal{O}(Q_2)$

Controlling Q

Controlling Q

Finally we have the following asymptotic behaviours for the above quantities:

Controlling Q

Finally we have the following asymptotic behaviours for the above quantities:

Asymptotics of Q

Controlling Q

Finally we have the following asymptotic behaviours for the above quantities:

Asymptotics of Q

Take fully connected network $f(\mathbf{w}; \mathbf{x})$ with initialization \mathbf{w}_0 , $\mathbf{w} \in \mathcal{B}(\mathbf{w}_0, R)$ with $R > 0$. Then with high probability over initialization:

Controlling \mathcal{Q}

Finally we have the following asymptotic behaviours for the above quantities:

Asymptotics of \mathcal{Q}

Take fully connected network $f(\mathbf{w}; \mathbf{x})$ with initialization \mathbf{w}_0 , $\mathbf{w} \in \mathcal{B}(\mathbf{w}_0, R)$ with $R > 0$. Then with high probability over initialization:

- $\mathcal{Q}_\infty = \tilde{O}\left(\frac{1}{\sqrt{m}}\right)$

Controlling \mathcal{Q}

Finally we have the following asymptotic behaviours for the above quantities:

Asymptotics of \mathcal{Q}

Take fully connected network $f(\mathbf{w}; \mathbf{x})$ with initialization \mathbf{w}_0 , $\mathbf{w} \in \mathcal{B}(\mathbf{w}_0, R)$ with $R > 0$. Then with high probability over initialization:

- $\mathcal{Q}_\infty = \tilde{\mathcal{O}}\left(\frac{1}{\sqrt{m}}\right)$
- $\mathcal{Q}_{2,2,1} = \mathcal{O}(1)$

Controlling \mathcal{Q}

Finally we have the following asymptotic behaviours for the above quantities:

Asymptotics of \mathcal{Q}

Take fully connected network $f(\mathbf{w}; \mathbf{x})$ with initialization \mathbf{w}_0 , $\mathbf{w} \in \mathcal{B}(\mathbf{w}_0, R)$ with $R > 0$. Then with high probability over initialization:

- $\mathcal{Q}_\infty = \tilde{\mathcal{O}}\left(\frac{1}{\sqrt{m}}\right)$
- $\mathcal{Q}_{2,2,1} = \mathcal{O}(1)$
- $\mathcal{Q}_2 = \mathcal{O}(1)$

Main Theorem

Main Theorem

Now we can put all together to obtain as a **direct** consequence:

Main Theorem

Now we can put all together to obtain as a **direct** consequence:

Asymptotics of Hessian

Main Theorem

Now we can put all together to obtain as a **direct** consequence:

Asymptotics of Hessian

Take fully connected network $f(\mathbf{w}; \mathbf{x})$ with initialization \mathbf{w}_0 , $\mathbf{w} \in \mathcal{B}(\mathbf{w}_0, R)$ with $R > 0$. Then with high probability over initialization:

Main Theorem

Now we can put all together to obtain as a **direct** consequence:

Asymptotics of Hessian

Take fully connected network $f(\mathbf{w}; \mathbf{x})$ with initialization \mathbf{w}_0 , $\mathbf{w} \in \mathcal{B}(\mathbf{w}_0, R)$ with $R > 0$. Then with high probability over initialization:

$$\|\mathbf{H}\|_2 = \tilde{O}\left(\frac{1}{\sqrt{m}}\right)$$

Main Theorem

Now we can put all together to obtain as a **direct** consequence:

Asymptotics of Hessian

Take fully connected network $f(\mathbf{w}; \mathbf{x})$ with initialization \mathbf{w}_0 , $\mathbf{w} \in \mathcal{B}(\mathbf{w}_0, R)$ with $R > 0$. Then with high probability over initialization:

$$\|\mathbf{H}\|_2 = \tilde{O}\left(\frac{1}{\sqrt{m}}\right)$$

- Fixing $R > 0$ big enough to include GD trajectory and applying the Lemma from before implies the **constancy** of the tangent kernel during **training**.

Main Theorem

Now we can put all together to obtain as a **direct** consequence:

Asymptotics of Hessian

Take fully connected network $f(\mathbf{w}; \mathbf{x})$ with initialization \mathbf{w}_0 , $\mathbf{w} \in \mathcal{B}(\mathbf{w}_0, R)$ with $R > 0$. Then with high probability over initialization:

$$\|\mathbf{H}\|_2 = \tilde{\mathcal{O}}\left(\frac{1}{\sqrt{m}}\right)$$

- Fixing $R > 0$ big enough to include GD trajectory and applying the Lemma from before implies the **constancy** of the tangent kernel during **training**.
- $\|\mathbf{K}\|_2 = \Theta(1)$ since $\mathcal{Q}_2 = \mathcal{O}(1)$

Some Comments

Some Comments

- The role of the radius R is a bit unclear.

Some Comments

- The role of the radius R is a bit unclear.
- What dependence on R in the bound for $\|\mathbf{H}\|_2$?

Some Comments

- The role of the radius R is a bit unclear.
- What dependence on R in the bound for $\|\mathbf{H}\|_2$?
- Probably all hidden in \mathcal{O} but for instance $R = \infty$ shouldn't work? Or only low probability bound?

Some Comments

- The role of the radius R is a bit unclear.
- What dependence on R in the bound for $\|\mathbf{H}\|_2$?
- Probably all hidden in \mathcal{O} but for instance $R = \infty$ shouldn't work? Or only low probability bound?
- How does the choice of R as GD radius affect the convergence/probability guarantee?

Negative Result for Lazy Learning

Negative Result for Lazy Learning

- Lots of literature claims that constancy of NTK is a consequence of the setting (**lazy learning**):

Negative Result for Lazy Learning

- Lots of literature claims that constancy of NTK is a consequence of the setting (**lazy learning**):

Small Learning rate + Huge Widths + Small Parameter Change

Negative Result for Lazy Learning

- Lots of literature claims that constancy of NTK is a consequence of the setting (**lazy learning**):

Small Learning rate + Huge Widths + Small Parameter Change

- Authors show that using a non-linear output layer instead of a linear layer induces a non-constant NTK

Negative Result for Lazy Learning

- Lots of literature claims that constancy of NTK is a consequence of the setting (**lazy learning**):

Small Learning rate + Huge Widths + Small Parameter Change

- Authors show that using a non-linear output layer instead of a linear layer induces a non-constant NTK
- Namely, using the 1-dim 1-hidden layer network $f(\mathbf{w}; x)$ as

$$\tilde{f}(\mathbf{w}, x) = \phi(f(\mathbf{w}; x))$$

where $\phi : \mathbb{R} \rightarrow \mathbb{R}$ is some twice-differentiable non-linearity

NTK is Not Constant (I)

NTK is Not Constant (I)

- Let's look at the Hessian:

NTK is Not Constant (I)

- Let's look at the Hessian:

$$\tilde{H}_{ij} = \frac{\partial^2 \tilde{f}}{\partial w_i \partial w_j}$$

NTK is Not Constant (I)

- Let's look at the Hessian:

$$\tilde{H}_{ij} = \frac{\partial^2 \tilde{f}}{\partial w_i \partial w_j} = \frac{\partial^2}{\partial w_i \partial w_j} \phi(f)$$

NTK is Not Constant (I)

- Let's look at the Hessian:

$$\tilde{H}_{ij} = \frac{\partial^2 \tilde{f}}{\partial w_i \partial w_j} = \frac{\partial^2}{\partial w_i \partial w_j} \phi(f) = \frac{\partial}{\partial w_i} \left(\phi'(f) \frac{\partial f}{\partial w_j} \right)$$

NTK is Not Constant (I)

- Let's look at the Hessian:

$$\begin{aligned}\tilde{H}_{ij} &= \frac{\partial^2 \tilde{f}}{\partial w_i \partial w_j} = \frac{\partial^2}{\partial w_i \partial w_j} \phi(f) = \frac{\partial}{\partial w_i} \left(\phi'(f) \frac{\partial f}{\partial w_j} \right) \\ &= \phi''(f) \frac{\partial f}{\partial w_i} \frac{\partial f}{\partial w_j} + \phi'(f) \frac{\partial^2 f}{\partial w_j \partial w_i}\end{aligned}$$

NTK is Not Constant (I)

- Let's look at the Hessian:

$$\begin{aligned}\tilde{H}_{ij} &= \frac{\partial^2 \tilde{f}}{\partial w_i \partial w_j} = \frac{\partial^2}{\partial w_i \partial w_j} \phi(f) = \frac{\partial}{\partial w_i} \left(\phi'(f) \frac{\partial f}{\partial w_j} \right) \\ &= \phi''(f) \frac{\partial f}{\partial w_i} \frac{\partial f}{\partial w_j} + \phi'(f) \frac{\partial^2 f}{\partial w_j \partial w_i}\end{aligned}$$

- Thus:**

NTK is Not Constant (I)

- Let's look at the Hessian:

$$\begin{aligned}\tilde{H}_{ij} &= \frac{\partial^2 \tilde{f}}{\partial w_i \partial w_j} = \frac{\partial^2}{\partial w_i \partial w_j} \phi(f) = \frac{\partial}{\partial w_i} \left(\phi'(f) \frac{\partial f}{\partial w_j} \right) \\ &= \phi''(f) \frac{\partial f}{\partial w_i} \frac{\partial f}{\partial w_j} + \phi'(f) \frac{\partial^2 f}{\partial w_j \partial w_i}\end{aligned}$$

- Thus:** $\tilde{\mathbf{H}} = \phi''(f) (\nabla_{\mathbf{w}} f) (\nabla_{\mathbf{w}} f)^T + \phi'(f) \mathbf{H}$

NTK is Not Constant (I)

- Let's look at the Hessian:

$$\begin{aligned}\tilde{H}_{ij} &= \frac{\partial^2 \tilde{f}}{\partial w_i \partial w_j} = \frac{\partial^2}{\partial w_i \partial w_j} \phi(f) = \frac{\partial}{\partial w_i} \left(\phi'(f) \frac{\partial f}{\partial w_j} \right) \\ &= \phi''(f) \frac{\partial f}{\partial w_i} \frac{\partial f}{\partial w_j} + \phi'(f) \frac{\partial^2 f}{\partial w_j \partial w_i}\end{aligned}$$

- Thus:** $\tilde{\mathbf{H}} = \phi''(f) (\nabla_{\mathbf{w}} f) (\nabla_{\mathbf{w}} f)^T + \phi'(f) \mathbf{H}$
- From **Weyl's**-inequality we have that for all symmetric matrices $\mathbf{A} = \mathbf{B} + \mathbf{C}$:

NTK is Not Constant (I)

- Let's look at the Hessian:

$$\begin{aligned}\tilde{H}_{ij} &= \frac{\partial^2 \tilde{f}}{\partial w_i \partial w_j} = \frac{\partial^2}{\partial w_i \partial w_j} \phi(f) = \frac{\partial}{\partial w_i} \left(\phi'(f) \frac{\partial f}{\partial w_j} \right) \\ &= \phi''(f) \frac{\partial f}{\partial w_i} \frac{\partial f}{\partial w_j} + \phi'(f) \frac{\partial^2 f}{\partial w_j \partial w_i}\end{aligned}$$

- Thus:** $\tilde{\mathbf{H}} = \phi''(f) (\nabla_{\mathbf{w}} f) (\nabla_{\mathbf{w}} f)^T + \phi'(f) \mathbf{H}$
- From **Weyl's**-inequality we have that for all symmetric matrices $\mathbf{A} = \mathbf{B} + \mathbf{C}$:

$$\lambda_1(\mathbf{A}) \geq \lambda_1(\mathbf{B}) + \lambda_n(\mathbf{C}) \iff \|\mathbf{A}\|_2 \geq \|\mathbf{B}\|_2 - \|\mathbf{C}\|_2$$

NTK is Not Constant (I)

- Let's look at the Hessian:

$$\begin{aligned}\tilde{H}_{ij} &= \frac{\partial^2 \tilde{f}}{\partial w_i \partial w_j} = \frac{\partial^2}{\partial w_i \partial w_j} \phi(f) = \frac{\partial}{\partial w_i} \left(\phi'(f) \frac{\partial f}{\partial w_j} \right) \\ &= \phi''(f) \frac{\partial f}{\partial w_i} \frac{\partial f}{\partial w_j} + \phi'(f) \frac{\partial^2 f}{\partial w_j \partial w_i}\end{aligned}$$

- Thus:** $\tilde{\mathbf{H}} = \phi''(f) (\nabla_{\mathbf{w}} f) (\nabla_{\mathbf{w}} f)^T + \phi'(f) \mathbf{H}$
- From **Weyl's**-inequality we have that for all symmetric matrices $\mathbf{A} = \mathbf{B} + \mathbf{C}$:

$$\lambda_1(\mathbf{A}) \geq \lambda_1(\mathbf{B}) + \lambda_n(\mathbf{C}) \iff \|\mathbf{A}\|_2 \geq \|\mathbf{B}\|_2 - \|\mathbf{C}\|_2$$

- Applying this to $\tilde{\mathbf{H}}$ leads to

NTK is Not Constant (I)

- Let's look at the Hessian:

$$\begin{aligned}\tilde{H}_{ij} &= \frac{\partial^2 \tilde{f}}{\partial w_i \partial w_j} = \frac{\partial^2}{\partial w_i \partial w_j} \phi(f) = \frac{\partial}{\partial w_i} \left(\phi'(f) \frac{\partial f}{\partial w_j} \right) \\ &= \phi''(f) \frac{\partial f}{\partial w_i} \frac{\partial f}{\partial w_j} + \phi'(f) \frac{\partial^2 f}{\partial w_j \partial w_i}\end{aligned}$$

- Thus:** $\tilde{\mathbf{H}} = \phi''(f) (\nabla_{\mathbf{w}} f) (\nabla_{\mathbf{w}} f)^T + \phi'(f) \mathbf{H}$
- From **Weyl's**-inequality we have that for all symmetric matrices $\mathbf{A} = \mathbf{B} + \mathbf{C}$:

$$\lambda_1(\mathbf{A}) \geq \lambda_1(\mathbf{B}) + \lambda_n(\mathbf{C}) \iff \|\mathbf{A}\|_2 \geq \|\mathbf{B}\|_2 - \|\mathbf{C}\|_2$$

- Applying this to $\tilde{\mathbf{H}}$ leads to

$$\|\tilde{\mathbf{H}}\|_2 \geq |\phi''(f)| \|\nabla_{\mathbf{w}} f\|_2^2 - |\phi'(f)| \|\mathbf{H}\|_2$$

NTK is Not Constant (II)

NTK is Not Constant (II)

- We found that

$$\|\tilde{\mathbf{H}}\|_2 \geq |\phi''(f)| \|\nabla_{\mathbf{w}} f\|_2^2 - |\phi'(f)| \|\mathbf{H}\|_2$$

NTK is Not Constant (II)

- We found that

$$\|\tilde{\mathbf{H}}\|_2 \geq |\phi''(f)| \|\nabla_{\mathbf{w}} f\|_2^2 - |\phi'(f)| \|\mathbf{H}\|_2$$

- Recall: $\|\mathbf{H}\|_2 \xrightarrow{m \rightarrow \infty} 0$

NTK is Not Constant (II)

- We found that

$$\|\tilde{\mathbf{H}}\|_2 \geq |\phi''(f)| \|\nabla_{\mathbf{w}} f\|_2^2 - |\phi'(f)| \|\mathbf{H}\|_2$$

- **Recall:** $\|\mathbf{H}\|_2 \xrightarrow{m \rightarrow \infty} 0$
- **However:** $\|\nabla_{\mathbf{w}} f\|_2^2 = K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \Theta(1)$

NTK is Not Constant (II)

- We found that

$$\|\tilde{\mathbf{H}}\|_2 \geq |\phi''(f)| \|\nabla_{\mathbf{w}} f\|_2^2 - |\phi'(f)| \|\mathbf{H}\|_2$$

- **Recall:** $\|\mathbf{H}\|_2 \xrightarrow{m \rightarrow \infty} 0$
- **However:** $\|\nabla_{\mathbf{w}} f\|_2^2 = K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}) = \Theta(1)$
- **Thus,** as long as ϕ'' doesn't vanish (so if not linear) we have

$$\|\tilde{\mathbf{H}}\|_2 = \Omega(1)$$

Other Example

Other Example

- Non-constancy of the NTK also for a **bottle-neck** architecture

Other Example

- Non-constancy of the NTK also for a **bottle-neck** architecture
- Specifically, they assume one hidden layer l does **not** go to infinity but $\alpha^{(l)}(\mathbf{x}) \in \mathbb{R}$, so $m_l = 1$ fixed and

$$\alpha^{(l+1)} = \sigma(\mathbf{w}^{(l+1)} \alpha^{(l)})$$

Other Example

- Non-constancy of the NTK also for a **bottle-neck** architecture
- Specifically, they assume one hidden layer l does **not** go to infinity but $\alpha^{(l)}(\mathbf{x}) \in \mathbb{R}$, so $m_l = 1$ fixed and

$$\alpha^{(l+1)} = \sigma(\mathbf{w}^{(l+1)} \alpha^{(l)})$$

- **But then** a term in $\mathcal{Q}_{2,2,1}$ scales as

=

Other Example

- Non-constancy of the NTK also for a **bottle-neck** architecture
- Specifically, they assume one hidden layer l does **not** go to infinity but $\alpha^{(l)}(\mathbf{x}) \in \mathbb{R}$, so $m_l = 1$ fixed and

$$\alpha^{(l+1)} = \sigma(\mathbf{w}^{(l+1)} \alpha^{(l)})$$

- **But then** a term in $\mathcal{Q}_{2,2,1}$ scales as

$$\left\| \frac{\partial \alpha^{(l+1)}}{(\partial \alpha^{(l)})^2} \right\|_{2,2,1} =$$

Other Example

- Non-constancy of the NTK also for a **bottle-neck** architecture
- Specifically, they assume one hidden layer l does **not** go to infinity but $\alpha^{(l)}(\mathbf{x}) \in \mathbb{R}$, so $m_l = 1$ fixed and

$$\alpha^{(l+1)} = \sigma(\mathbf{w}^{(l+1)} \alpha^{(l)})$$

- **But then** a term in $\mathcal{Q}_{2,2,1}$ scales as

$$\left\| \frac{\partial \alpha^{(l+1)}}{(\partial \alpha^{(l)})^2} \right\|_{2,2,1} = \sum_{i=1}^m \left| \sigma'' \left(w_i^{(l+1)} \alpha^{(l)} \right) \left(w_i^{(l+1)} \right)^2 \right|$$

Other Example

- Non-constancy of the NTK also for a **bottle-neck** architecture
- Specifically, they assume one hidden layer l does **not** go to infinity but $\alpha^{(l)}(\mathbf{x}) \in \mathbb{R}$, so $m_l = 1$ fixed and

$$\alpha^{(l+1)} = \sigma(\mathbf{w}^{(l+1)} \alpha^{(l)})$$

- **But then** a term in $\mathcal{Q}_{2,2,1}$ scales as

$$\left\| \frac{\partial \alpha^{(l+1)}}{(\partial \alpha^{(l)})^2} \right\|_{2,2,1} = \sum_{i=1}^m \left| \sigma'' \left(w_i^{(l+1)} \alpha^{(l)} \right) \left(w_i^{(l+1)} \right)^2 \right| = \Theta(m)$$

Other Example

- Non-constancy of the NTK also for a **bottle-neck** architecture
- Specifically, they assume one hidden layer l does **not** go to infinity but $\alpha^{(l)}(\mathbf{x}) \in \mathbb{R}$, so $m_l = 1$ fixed and

$$\alpha^{(l+1)} = \sigma(\mathbf{w}^{(l+1)} \alpha^{(l)})$$

- **But then** a term in $\mathcal{Q}_{2,2,1}$ scales as

$$\left\| \frac{\partial \alpha^{(l+1)}}{(\partial \alpha^{(l)})^2} \right\|_{2,2,1} = \sum_{i=1}^m \left| \sigma'' \left(w_i^{(l+1)} \alpha^{(l)} \right) \left(w_i^{(l+1)} \right)^2 \right| = \Theta(m)$$

- This "kills" the $\frac{1}{\sqrt{m}}$ from before!

Overview

Overview

- In summary we found that

Overview

- In summary we found that

Model	$\left\ \frac{\partial f}{\partial \alpha^{(l)}} \right\ _{\infty}$	(2, 2, 1)-norms	Hessian norm	Trans. to linearity?
linear output layer	$\tilde{O}(1/\sqrt{m})$	$O(1)$	$\tilde{O}(1/\sqrt{m})$	Yes
non-linear output layer	$\tilde{O}(1)$	$O(1)$	$\tilde{O}(1)$	No
bottleneck	$\tilde{O}(1/\sqrt{m})$	$O(m)$	$\tilde{O}(1)$	No

Overview

- In summary we found that

Model	$\left\ \frac{\partial f}{\partial \alpha^{(l)}} \right\ _{\infty}$	(2, 2, 1)-norms	Hessian norm	Trans. to linearity?
linear output layer	$\tilde{O}(1/\sqrt{m})$	$O(1)$	$\tilde{O}(1/\sqrt{m})$	Yes
non-linear output layer	$\tilde{O}(1)$	$O(1)$	$\tilde{O}(1)$	No
bottleneck	$\tilde{O}(1/\sqrt{m})$	$O(m)$	$\tilde{O}(1)$	No

- One important side note. For GD to converge fast, it is not necessary that the tangent kernel becomes constant but that it remains **well-conditioned** along the optimization path.

Overview

- In summary we found that

Model	$\left\ \frac{\partial f}{\partial \alpha^{(l)}} \right\ _{\infty}$	(2, 2, 1)-norms	Hessian norm	Trans. to linearity?
linear output layer	$\tilde{O}(1/\sqrt{m})$	$O(1)$	$\tilde{O}(1/\sqrt{m})$	Yes
non-linear output layer	$\tilde{O}(1)$	$O(1)$	$\tilde{O}(1)$	No
bottleneck	$\tilde{O}(1/\sqrt{m})$	$O(m)$	$\tilde{O}(1)$	No

- One important side note. For GD to converge fast, it is not necessary that the tangent kernel becomes constant but that it remains **well-conditioned** along the optimization path.
- Well-conditionedness is somewhat **"passed"** on to the non-linear output model from the linear output model. That's why these models also converge fast in practice.

Experiments (I)

Experiments (I)

- Track the change of the tangent kernel by computing

Experiments (I)

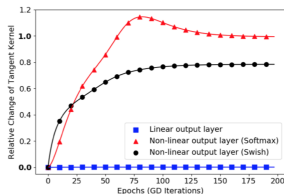
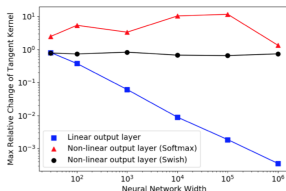
- Track the change of the tangent kernel by computing

$$\Delta K = \sup_{t>0} \frac{\|K(\mathbf{w}_t) - K(\mathbf{w}_0)\|_F}{\|K(\mathbf{w}_0)\|_F}$$

Experiments (I)

- Track the change of the tangent kernel by computing

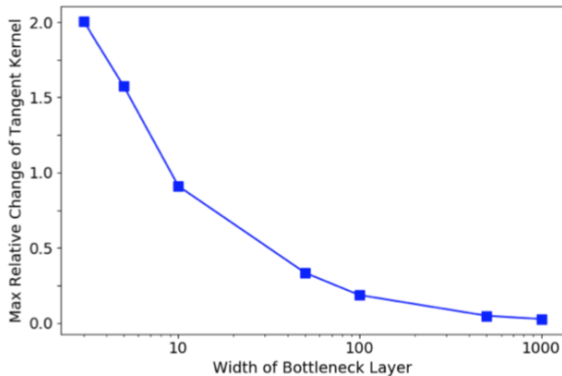
$$\Delta K = \sup_{t>0} \frac{\|K(\mathbf{w}_t) - K(\mathbf{w}_0)\|_F}{\|K(\mathbf{w}_0)\|_F}$$



Change in NTK for **linear** versus **non-Linear** output activation. 1 hidden layer network with either linear, softmax or swish activation at output.

Experiments (II)

Experiments (II)



Width of the bottleneck layer and corresponding change in NTK.
3 hidden layer network of widths $m = 10^4$

Closing Thoughts

Closing Thoughts

- Overall very interesting and intuitive formulation

Closing Thoughts

- Overall very interesting and intuitive formulation
- Very elementary maths (Zero Hessian \implies Linear etc)

Closing Thoughts

- Overall very interesting and intuitive formulation
- Very elementary maths (Zero Hessian \implies Linear etc)
- Of course lots of math hidden in the statement that GD only explores $\mathcal{B}(\mathbf{w}_0, R)$. Some mistakes in the proofs but I think only little fixes needed.

Closing Thoughts

- Overall very interesting and intuitive formulation
- Very elementary maths (Zero Hessian \implies Linear etc)
- Of course lots of math hidden in the statement that GD only explores $\mathcal{B}(\mathbf{w}_0, R)$. Some mistakes in the proofs but I think only little fixes needed.
- Very interesting that the optimization mechanism is clearly **separated**. I have not understood the paper in detail but seems like it suffices to show (for instance for SGD) that the iterates stay in $\mathcal{B}(\mathbf{w}_0, R)$ for R independent of width m .

Closing Thoughts

- Overall very interesting and intuitive formulation
- Very elementary maths (Zero Hessian \implies Linear etc)
- Of course lots of math hidden in the statement that GD only explores $\mathcal{B}(\mathbf{w}_0, R)$. Some mistakes in the proofs but I think only little fixes needed.
- Very interesting that the optimization mechanism is clearly **separated**. I have not understood the paper in detail but seems like it suffices to show (for instance for SGD) that the iterates stay in $\mathcal{B}(\mathbf{w}_0, R)$ for R independent of width m .
- Are similar results known for SGD, ADAM, second order methods? What does it take to fall out of this ball (except for big learning rates)?