

# **Neural Tangent Kernel: Convergence and Generalization in Neural Networks**

**Arthur Jacot, Franck Gabriel, Clément Hongler**

# First Approach

- I will try to give **high-level** overview first

# First Approach

- I will try to give **high-level** overview first
- Different notation to not enter into all mathematical details  
*(dual space, functional kernel regression etc)*

# First Approach

- I will try to give **high-level** overview first
- Different notation to not enter into all mathematical details (*dual space, functional kernel regression etc*)
- If time allows we can dwelve deeper into the matter

# **Setup**

# Setup

- Data points  $(\mathbf{x}_i, y_i) \in \mathbb{R}^{d_0} \times \mathbb{R}$  for  $i = 1, \dots, n$

# Setup

- Data points  $(\mathbf{x}_i, y_i) \in \mathbb{R}^{d_0} \times \mathbb{R}$  for  $i = 1, \dots, n$
- Stack features and targets into matrices:

$$\mathbf{X} \in \mathbb{R}^{n \times d_0} \text{ and } \mathbf{y} \in \mathbb{R}^n$$

# Notation

# Notation

- Only consider **fully-connected** networks of depth  $L$ :

# Notation

- Only consider **fully-connected** networks of depth  $L$ :

$$\begin{aligned} - \tilde{\alpha}^{(l+1)}(\mathbf{x}) &= \frac{1}{\sqrt{d_l}} \mathbf{W}^{(l)} \alpha^{(l)}(\mathbf{x}) + \beta \mathbf{b}^{(l)} \in \mathbb{R}^{d_{l+1}} \\ - \alpha^{(l+1)}(\mathbf{x}) &= \sigma(\tilde{\alpha}^{(l+1)}(\mathbf{x})) \in \mathbb{R}^{d_{l+1}} \end{aligned}$$

# Notation

- Only consider **fully-connected** networks of depth  $L$ :

$$\begin{aligned} - \tilde{\alpha}^{(l+1)}(\mathbf{x}) &= \frac{1}{\sqrt{d_l}} \mathbf{W}^{(l)} \alpha^{(l)}(\mathbf{x}) + \beta \mathbf{b}^{(l)} \in \mathbb{R}^{d_{l+1}} \\ - \alpha^{(l+1)}(\mathbf{x}) &= \sigma(\tilde{\alpha}^{(l+1)}(\mathbf{x})) \in \mathbb{R}^{d_{l+1}} \end{aligned}$$

where  $\mathbf{W}^{(l)} \in \mathbb{R}^{d_{l+1} \times d_l}$ ,  $\mathbf{b}^{(l)} \in \mathbb{R}^{d_{l+1}}$ ,  $\sigma$  is a coordinate-wise activation function and  $\beta > 0$

# Notation

- Only consider **fully-connected** networks of depth  $L$ :

$$\begin{aligned} - \tilde{\alpha}^{(l+1)}(\mathbf{x}) &= \frac{1}{\sqrt{d_l}} \mathbf{W}^{(l)} \alpha^{(l)}(\mathbf{x}) + \beta \mathbf{b}^{(l)} \in \mathbb{R}^{d_{l+1}} \\ - \alpha^{(l+1)}(\mathbf{x}) &= \sigma(\tilde{\alpha}^{(l+1)}(\mathbf{x})) \in \mathbb{R}^{d_{l+1}} \end{aligned}$$

where  $\mathbf{W}^{(l)} \in \mathbb{R}^{d_{l+1} \times d_l}$ ,  $\mathbf{b}^{(l)} \in \mathbb{R}^{d_{l+1}}$ ,  $\sigma$  is a coordinate-wise activation function and  $\beta > 0$

- $f_{\theta}(\mathbf{x}) = \tilde{\alpha}^{(L)}(\mathbf{x}) \in \mathbb{R}^{d_L}$  is the **output** of the network

# Notation

- Only consider **fully-connected** networks of depth  $L$ :

$$\begin{aligned} - \tilde{\alpha}^{(l+1)}(\mathbf{x}) &= \frac{1}{\sqrt{d_l}} \mathbf{W}^{(l)} \alpha^{(l)}(\mathbf{x}) + \beta \mathbf{b}^{(l)} \in \mathbb{R}^{d_{l+1}} \\ - \alpha^{(l+1)}(\mathbf{x}) &= \sigma(\tilde{\alpha}^{(l+1)}(\mathbf{x})) \in \mathbb{R}^{d_{l+1}} \end{aligned}$$

where  $\mathbf{W}^{(l)} \in \mathbb{R}^{d_{l+1} \times d_l}$ ,  $\mathbf{b}^{(l)} \in \mathbb{R}^{d_{l+1}}$ ,  $\sigma$  is a coordinate-wise activation function and  $\beta > 0$

- $f_{\theta}(\mathbf{x}) = \tilde{\alpha}^{(L)}(\mathbf{x}) \in \mathbb{R}^{d_L}$  is the **output** of the network
- Number of parameters  $P = \sum_{l=1}^L d_{l-1} d_l$

# Notation

- Only consider **fully-connected** networks of depth  $L$ :

$$\begin{aligned} - \tilde{\alpha}^{(l+1)}(\mathbf{x}) &= \frac{1}{\sqrt{d_l}} \mathbf{W}^{(l)} \alpha^{(l)}(\mathbf{x}) + \beta \mathbf{b}^{(l)} \in \mathbb{R}^{d_{l+1}} \\ - \alpha^{(l+1)}(\mathbf{x}) &= \sigma(\tilde{\alpha}^{(l+1)}(\mathbf{x})) \in \mathbb{R}^{d_{l+1}} \end{aligned}$$

where  $\mathbf{W}^{(l)} \in \mathbb{R}^{d_{l+1} \times d_l}$ ,  $\mathbf{b}^{(l)} \in \mathbb{R}^{d_{l+1}}$ ,  $\sigma$  is a coordinate-wise activation function and  $\beta > 0$

- $f_{\theta}(\mathbf{x}) = \tilde{\alpha}^{(L)}(\mathbf{x}) \in \mathbb{R}^{d_L}$  is the **output** of the network
- Number of parameters  $P = \sum_{l=1}^L d_{l-1} d_l$
- Same** notation as last time :)

# **Random Initialization**

# Random Initialization

- Initialize as  $W_{ij}^{(l)} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$  and  $b_i^{(l)} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$

# Random Initialization

- Initialize as  $W_{ij}^{(l)} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$  and  $b_i^{(l)} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$
- Made **width dependence** of random initialization explicit in the architecture

# Random Initialization

- Initialize as  $W_{ij}^{(l)} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$  and  $b_i^{(l)} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$
- Made **width dependence** of random initialization explicit in the architecture
- **Effectively:**  $W_{ij}^{(l)} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \frac{1}{d_l})$  and  $b_i^{(l)} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \beta^2)$   
 $\implies$  Usual Lecun initialization

# Random Initialization

- Initialize as  $W_{ij}^{(l)} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$  and  $b_i^{(l)} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$
- Made **width dependence** of random initialization explicit in the architecture
- **Effectively:**  $W_{ij}^{(l)} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \frac{1}{d_l})$  and  $b_i^{(l)} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \beta^2)$   
     $\implies$  Usual Lecun initialization
- Again **identical** to first presentation

# **Loss**

# Loss

- This time we will **train** the neural network

# Loss

- This time we will **train** the neural network
- We take some general **loss** function

$$L : \mathbb{R}^P \times \mathbb{R}^{n \times d_0} \times \mathbb{R}^n \rightarrow \mathbb{R}, \quad (\theta, \mathbf{X}, \mathbf{y}) \mapsto L(f_\theta(\mathbf{X}), \mathbf{y})$$

# Loss

- This time we will **train** the neural network
- We take some general **loss** function

$$L : \mathbb{R}^P \times \mathbb{R}^{n \times d_0} \times \mathbb{R}^n \rightarrow \mathbb{R}, \quad (\theta, \mathbf{X}, \mathbf{y}) \mapsto L(f_\theta(\mathbf{X}), \mathbf{y})$$

- Assume that the loss decomposes over the samples:

$$L(f_\theta(\mathbf{X}), \mathbf{y}) = \sum_{i=1}^n L(f_\theta(\mathbf{x}_i), y_i)$$

# Gradient Flow (I)

# Gradient Flow (I)

- We usually do **gradient descent**

$$\theta_{k+1} = \theta_k - \eta \nabla_{\theta} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta_k}$$

# Gradient Flow (I)

- We usually do **gradient descent**

$$\theta_{k+1} = \theta_k - \eta \nabla_{\theta} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta_k}$$

- Let's take a continuous view point:  $\theta_k$  are samples at multiples of  $\eta$  from a function  $\theta : \mathbb{R}_+ \rightarrow \mathbb{R}^{d_0}$

$$\theta_k = \theta(k\eta)$$

# Gradient Flow (I)

- We usually do **gradient descent**

$$\theta_{k+1} = \theta_k - \eta \nabla_{\theta} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta_k}$$

- Let's take a continuous view point:  $\theta_k$  are samples at multiples of  $\eta$  from a function  $\theta : \mathbb{R}_+ \rightarrow \mathbb{R}^{d_0}$

$$\theta_k = \theta(k\eta)$$

- Let  $t = k\eta$  and let's look at

$$\theta(t + \eta)$$

# Gradient Flow (I)

- We usually do **gradient descent**

$$\theta_{k+1} = \theta_k - \eta \nabla_{\theta} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta_k}$$

- Let's take a continuous view point:  $\theta_k$  are samples at multiples of  $\eta$  from a function  $\theta : \mathbb{R}_+ \rightarrow \mathbb{R}^{d_0}$

$$\theta_k = \theta(k\eta)$$

- Let  $t = k\eta$  and let's look at

$$\theta(t + \eta) = \theta_{k+1}$$

# Gradient Flow (I)

- We usually do **gradient descent**

$$\theta_{k+1} = \theta_k - \eta \nabla_{\theta} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta_k}$$

- Let's take a continuous view point:  $\theta_k$  are samples at multiples of  $\eta$  from a function  $\theta : \mathbb{R}_+ \rightarrow \mathbb{R}^{d_0}$

$$\theta_k = \theta(k\eta)$$

- Let  $t = k\eta$  and let's look at

$$\theta(t + \eta) = \theta_{k+1} = \theta_k - \eta \nabla_{\theta} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta_k}$$

# Gradient Flow (I)

- We usually do **gradient descent**

$$\theta_{k+1} = \theta_k - \eta \nabla_{\theta} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta_k}$$

- Let's take a continuous view point:  $\theta_k$  are samples at multiples of  $\eta$  from a function  $\theta : \mathbb{R}_+ \rightarrow \mathbb{R}^{d_0}$

$$\theta_k = \theta(k\eta)$$

- Let  $t = k\eta$  and let's look at

$$\begin{aligned}\theta(t + \eta) &= \theta_{k+1} = \theta_k - \eta \nabla_{\theta} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta_k} \\ &= \theta(t) - \eta \nabla_{\theta} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta(t)}\end{aligned}$$

# Gradient Flow (II)

## Gradient Flow (II)

- Rearranging gives

$$\frac{1}{\eta} (\theta(t + \eta) - \theta(t)) = -\nabla_{\theta} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta(t)}$$

## Gradient Flow (II)

- Rearranging gives

$$\frac{1}{\eta} (\theta(t + \eta) - \theta(t)) = -\nabla_{\theta} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta(t)}$$

- Letting  $\eta \rightarrow 0$  thus gives the gradient flow equation:

$$\dot{\theta}(t) = -\nabla_{\theta} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta(t)}$$

## Gradient Flow (II)

- Rearranging gives

$$\frac{1}{\eta} (\theta(t + \eta) - \theta(t)) = -\nabla_{\theta} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta(t)}$$

- Letting  $\eta \rightarrow 0$  thus gives the gradient flow equation:

$$\dot{\theta}(t) = -\nabla_{\theta} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta(t)}$$

- **Side note:** Gradient Descent is the Euler discretization of the above differential equation

## Gradient Flow (II)

- Rearranging gives

$$\frac{1}{\eta} (\theta(t + \eta) - \theta(t)) = -\nabla_{\theta} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta(t)}$$

- Letting  $\eta \rightarrow 0$  thus gives the gradient flow equation:

$$\dot{\theta}(t) = -\nabla_{\theta} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta(t)}$$

- **Side note:** Gradient Descent is the Euler discretization of the above differential equation
- **Non-linear ODE**, hard to solve... If we could solve it, no need to train NNs since the solution describes the trajectory through time

# Neural Tangent Kernel

Let us now introduce the quantity of interest:

**Definition 1:**

# Neural Tangent Kernel

Let us now introduce the quantity of interest:

## Definition 1:

The **Neural Tangent Kernel** is given by

$$\Theta_t^{(L)} = \sum_{p=1}^P \partial_{\theta_p} f_{\theta} \otimes \partial_{\theta_p} f_{\theta}$$

# Neural Tangent Kernel

Let us now introduce the quantity of interest:

**Definition 1:**

The **Neural Tangent Kernel** is given by

$$\Theta_t^{(L)} = \sum_{p=1}^P \partial_{\theta_p} f_{\theta} \otimes \partial_{\theta_p} f_{\theta}$$

This means in particular that for  $k, k' \in \{1, \dots, d_L\}$  and  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{d_0}$

$$\Theta_{kk'}^{(L)}(\mathbf{x}, \mathbf{x}') = \sum_{p=1}^P \partial_{\theta_p} f_{\theta, k}(\mathbf{x}) \partial_{\theta_p} f_{\theta, k'}(\mathbf{x}') = \left\langle \partial_{\theta} f_{\theta, k}(\mathbf{x}), \partial_{\theta} f_{\theta, k'}(\mathbf{x}') \right\rangle$$

# Neural Tangent Kernel

Let us now introduce the quantity of interest:

**Definition 1:**

The **Neural Tangent Kernel** is given by

$$\Theta_t^{(L)} = \sum_{p=1}^P \partial_{\theta_p} f_{\theta} \otimes \partial_{\theta_p} f_{\theta}$$

This means in particular that for  $k, k' \in \{1, \dots, d_L\}$  and  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{d_0}$

$$\Theta_{kk'}^{(L)}(\mathbf{x}, \mathbf{x}') = \sum_{p=1}^P \partial_{\theta_p} f_{\theta, k}(\mathbf{x}) \partial_{\theta_p} f_{\theta, k'}(\mathbf{x}') = \left\langle \partial_{\theta} f_{\theta, k}(\mathbf{x}), \partial_{\theta} f_{\theta, k'}(\mathbf{x}') \right\rangle$$

and thus  $\Theta^{(L)}(\mathbf{x}, \mathbf{x}') \in \mathbb{R}^{d_L \times d_L}$  (multi-dimensional kernel)

# An Easier Setting

# An Easier Setting

- Let's assume that  $d_L = 1$  (e.g. binary classification problem or least squared regression)

# An Easier Setting

- Let's assume that  $d_L = 1$  (e.g. binary classification problem or least squared regression)
- Then the NTK simplifies to a scalar:

$$\Theta_t^{(L)}(\mathbf{x}, \mathbf{x}') = \sum_{p=1}^P \partial_{\theta_p} f_{\theta}(\mathbf{x}) \partial_{\theta_p} f_{\theta}(\mathbf{x}')$$

# An Easier Setting

- Let's assume that  $d_L = 1$  (e.g. binary classification problem or least squared regression)
- Then the NTK simplifies to a scalar:

$$\Theta_t^{(L)}(\mathbf{x}, \mathbf{x}') = \sum_{p=1}^P \partial_{\theta_p} f_{\theta}(\mathbf{x}) \partial_{\theta_p} f_{\theta}(\mathbf{x}')$$

- $[\Theta_t^{(L)}(\mathbf{x}, \mathbf{X})]_i = \Theta_t^{(L)}(\mathbf{x}, \mathbf{x}_i)$  and  $[\Theta_t^{(L)}(\mathbf{X}, \mathbf{X})]_{ij} = \Theta_t^{(L)}(\mathbf{x}_i, \mathbf{x}_j)$

# An Easier Setting

- Let's assume that  $d_L = 1$  (e.g. binary classification problem or least squared regression)
- Then the NTK simplifies to a scalar:

$$\Theta_t^{(L)}(\mathbf{x}, \mathbf{x}') = \sum_{p=1}^P \partial_{\theta_p} f_{\theta}(\mathbf{x}) \partial_{\theta_p} f_{\theta}(\mathbf{x}')$$

- $[\Theta_t^{(L)}(\mathbf{x}, \mathbf{X})]_i = \Theta_t^{(L)}(\mathbf{x}, \mathbf{x}_i)$  and  $[\Theta_t^{(L)}(\mathbf{X}, \mathbf{X})]_{ij} = \Theta_t^{(L)}(\mathbf{x}_i, \mathbf{x}_j)$
- **Question:** Why should we care about this object???

# NTK and Gradient Flow

# NTK and Gradient Flow

- Let us "train"  $f_\theta$  with gradient flow:

# NTK and Gradient Flow

- Let us "train"  $f_\theta$  with gradient flow:

$$\dot{\theta}(t) = -\nabla_{\theta} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta(t)}$$

# NTK and Gradient Flow

- Let us "train"  $f_\theta$  with gradient flow:

$$\dot{\theta}(t) = -\nabla_{\theta} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta(t)}$$

- Component-wise:  $\frac{\partial \theta_p}{\partial t} = -\frac{\partial}{\partial \theta_p} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta(t)}$

# NTK and Gradient Flow

- Let us "train"  $f_\theta$  with gradient flow:

$$\dot{\theta}(t) = -\nabla_{\theta} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta(t)}$$

- Component-wise:  $\frac{\partial \theta_p}{\partial t} = -\frac{\partial}{\partial \theta_p} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta(t)}$
- This equation describes the dynamics in **weight space**

# NTK and Gradient Flow

- Let us "train"  $f_\theta$  with gradient flow:

$$\dot{\theta}(t) = -\nabla_{\theta} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta(t)}$$

- Component-wise:  $\frac{\partial \theta_p}{\partial t} = -\frac{\partial}{\partial \theta_p} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \Big|_{\theta=\theta(t)}$
- This equation describes the dynamics in **weight space**
- What can we say about the **function space** dynamics?

# Function Space Dynamics (I)

# Function Space Dynamics (I)

Let's take a new input  $x \in \mathbb{R}^{d_0}$ :

# Function Space Dynamics (I)

Let's take a new input  $\mathbf{x} \in \mathbb{R}^{d_0}$ :

$$\partial_t f_{\theta}(\mathbf{x})$$

# Function Space Dynamics (I)

Let's take a new input  $\mathbf{x} \in \mathbb{R}^{d_0}$ :

$$\partial_t f_{\theta}(\mathbf{x}) = \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \frac{\partial \theta_p}{\partial t}$$

# Function Space Dynamics (I)

Let's take a new input  $\mathbf{x} \in \mathbb{R}^{d_0}$ :

$$\partial_t f_{\theta}(\mathbf{x}) = \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \frac{\partial \theta_p}{\partial t} = - \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \partial_{\theta_p} L(f_{\theta}(\mathbf{X}), \mathbf{y})$$

# Function Space Dynamics (I)

Let's take a new input  $\mathbf{x} \in \mathbb{R}^{d_0}$ :

$$\begin{aligned}\partial_t f_{\theta}(\mathbf{x}) &= \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \frac{\partial \theta_p}{\partial t} = - \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \partial_{\theta_p} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \\ &= - \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \partial_{\theta_p} \sum_{i=1}^n L(f_{\theta}(\mathbf{x}_i), y_i)\end{aligned}$$

# Function Space Dynamics (I)

Let's take a new input  $\mathbf{x} \in \mathbb{R}^{d_0}$ :

$$\begin{aligned}\partial_t f_{\theta}(\mathbf{x}) &= \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \frac{\partial \theta_p}{\partial t} = - \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \partial_{\theta_p} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \\ &= - \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \partial_{\theta_p} \sum_{i=1}^n L(f_{\theta}(\mathbf{x}_i), y_i) \\ &= - \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \sum_{i=1}^n \frac{\partial L(f_{\theta}(\mathbf{x}_i), y_i)}{\partial f_{\theta}(\mathbf{x}_i)} \frac{\partial f_{\theta}(\mathbf{x}_i)}{\partial \theta_p}\end{aligned}$$

# Function Space Dynamics (I)

Let's take a new input  $\mathbf{x} \in \mathbb{R}^{d_0}$ :

$$\begin{aligned}\partial_t f_{\theta}(\mathbf{x}) &= \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \frac{\partial \theta_p}{\partial t} = - \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \partial_{\theta_p} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \\ &= - \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \partial_{\theta_p} \sum_{i=1}^n L(f_{\theta}(\mathbf{x}_i), y_i) \\ &= - \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \sum_{i=1}^n \frac{\partial L(f_{\theta}(\mathbf{x}_i), y_i)}{\partial f_{\theta}(\mathbf{x}_i)} \frac{\partial f_{\theta}(\mathbf{x}_i)}{\partial \theta_p} \\ &= - \sum_{i=1}^n \frac{\partial L(f_{\theta}(\mathbf{x}_i), y_i)}{\partial f_{\theta}(\mathbf{x}_i)} \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \frac{\partial f_{\theta}(\mathbf{x}_i)}{\partial \theta_p}\end{aligned}$$

# Function Space Dynamics (I)

Let's take a new input  $\mathbf{x} \in \mathbb{R}^{d_0}$ :

$$\begin{aligned}\partial_t f_{\theta}(\mathbf{x}) &= \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \frac{\partial \theta_p}{\partial t} = - \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \partial_{\theta_p} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \\ &= - \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \partial_{\theta_p} \sum_{i=1}^n L(f_{\theta}(\mathbf{x}_i), y_i) \\ &= - \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \sum_{i=1}^n \frac{\partial L(f_{\theta}(\mathbf{x}_i), y_i)}{\partial f_{\theta}(\mathbf{x}_i)} \frac{\partial f_{\theta}(\mathbf{x}_i)}{\partial \theta_p} \\ &= - \sum_{i=1}^n \frac{\partial L(f_{\theta}(\mathbf{x}_i), y_i)}{\partial f_{\theta}(\mathbf{x}_i)} \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \frac{\partial f_{\theta}(\mathbf{x}_i)}{\partial \theta_p} \\ &= - \sum_{i=1}^n \frac{\partial L(f_{\theta}(\mathbf{x}_i), y_i)}{\partial f_{\theta}(\mathbf{x}_i)} \Theta_t^{(L)}(\mathbf{x}, \mathbf{x}_i)\end{aligned}$$

# Function Space Dynamics (I)

Let's take a new input  $\mathbf{x} \in \mathbb{R}^{d_0}$ :

$$\begin{aligned}\partial_t f_{\theta}(\mathbf{x}) &= \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \frac{\partial \theta_p}{\partial t} = - \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \partial_{\theta_p} L(f_{\theta}(\mathbf{X}), \mathbf{y}) \\ &= - \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \partial_{\theta_p} \sum_{i=1}^n L(f_{\theta}(\mathbf{x}_i), y_i) \\ &= - \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \sum_{i=1}^n \frac{\partial L(f_{\theta}(\mathbf{x}_i), y_i)}{\partial f_{\theta}(\mathbf{x}_i)} \frac{\partial f_{\theta}(\mathbf{x}_i)}{\partial \theta_p} \\ &= - \sum_{i=1}^n \frac{\partial L(f_{\theta}(\mathbf{x}_i), y_i)}{\partial f_{\theta}(\mathbf{x}_i)} \sum_{p=1}^P \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta_p} \frac{\partial f_{\theta}(\mathbf{x}_i)}{\partial \theta_p} \\ &= - \sum_{i=1}^n \frac{\partial L(f_{\theta}(\mathbf{x}_i), y_i)}{\partial f_{\theta}(\mathbf{x}_i)} \Theta_t^{(L)}(\mathbf{x}, \mathbf{x}_i) = - \Theta_t^{(L)}(\mathbf{x}, \mathbf{X})^T \frac{\partial L(f_{\theta}(\mathbf{X}), \mathbf{y})}{\partial f_{\theta}(\mathbf{X})}\end{aligned}$$

# Function Space Dynamics (II)

## Function Space Dynamics (II)

- The function space dynamics is hence governed by the NTK:

## Function Space Dynamics (II)

- The function space dynamics is hence governed by the NTK:

$$\partial_t f_{\theta}(\mathbf{x}) = -\Theta_t^{(L)}(\mathbf{x}, \mathbf{X})^T \frac{\partial L(f_{\theta}(\mathbf{X}), \mathbf{y})}{\partial f_{\theta}(\mathbf{X})}$$

## Function Space Dynamics (II)

- The function space dynamics is hence governed by the NTK:

$$\partial_t f_{\theta}(\mathbf{x}) = -\Theta_t^{(L)}(\mathbf{x}, \mathbf{X})^T \frac{\partial L(f_{\theta}(\mathbf{X}), \mathbf{y})}{\partial f_{\theta}(\mathbf{X})}$$

- Notice:** This result so far holds for **any** neural network and **any** loss  $L$ , no infinite-width limit taken yet!

## Function Space Dynamics (II)

- The function space dynamics is hence governed by the NTK:

$$\partial_t f_{\theta}(\mathbf{x}) = -\Theta_t^{(L)}(\mathbf{x}, \mathbf{X})^T \frac{\partial L(f_{\theta}(\mathbf{X}), \mathbf{y})}{\partial f_{\theta}(\mathbf{X})}$$

- Notice:** This result so far holds for **any** neural network and **any** loss  $L$ , no infinite-width limit taken yet!
- $\Theta_t^{(L)}(\mathbf{x}, \mathbf{x}')$  is a **random variable** (Depends on initialization)

## Function Space Dynamics (II)

- The function space dynamics is hence governed by the NTK:

$$\partial_t f_{\theta}(\mathbf{x}) = -\Theta_t^{(L)}(\mathbf{x}, \mathbf{X})^T \frac{\partial L(f_{\theta}(\mathbf{X}), \mathbf{y})}{\partial f_{\theta}(\mathbf{X})}$$

- Notice:** This result so far holds for **any** neural network and **any** loss  $L$ , no infinite-width limit taken yet!
- $\Theta_t^{(L)}(\mathbf{x}, \mathbf{x}')$  is a **random variable** (Depends on initialization)
- $\Theta_t^{(L)}(\mathbf{x}, \mathbf{x}')$  is **time-dependent** (It is a function of the parameters  $\theta(t)$ )

## Function Space Dynamics (II)

- The function space dynamics is hence governed by the NTK:

$$\partial_t f_{\theta}(\mathbf{x}) = -\Theta_t^{(L)}(\mathbf{x}, \mathbf{X})^T \frac{\partial L(f_{\theta}(\mathbf{X}), \mathbf{y})}{\partial f_{\theta}(\mathbf{X})}$$

- Notice:** This result so far holds for **any** neural network and **any** loss  $L$ , no infinite-width limit taken yet!
- $\Theta_t^{(L)}(\mathbf{x}, \mathbf{x}')$  is a **random variable** (Depends on initialization)
- $\Theta_t^{(L)}(\mathbf{x}, \mathbf{x}')$  is **time-dependent** (It is a function of the parameters  $\theta(t)$ )
- ODE hence **not** analytically solvable (again if it were, no need to train NNs at all)

# Infinite-Width Limit

# Infinite-Width Limit

**Question:** What happens if we let  $d_1, \dots, d_{L-1} \rightarrow \infty$ ?

# Infinite-Width Limit

**Question:** What happens if we let  $d_1, \dots, d_{L-1} \rightarrow \infty$ ?

We recover the following properties:

# Infinite-Width Limit

**Question:** What happens if we let  $d_1, \dots, d_{L-1} \rightarrow \infty$ ?

We recover the following properties:

- $\Theta_t^{(L)}$  converges to a **deterministic** kernel  $\Theta_{\infty,t}^{(L)}$

# Infinite-Width Limit

**Question:** What happens if we let  $d_1, \dots, d_{L-1} \rightarrow \infty$ ?

We recover the following properties:

- $\Theta_t^{(L)}$  converges to a **deterministic** kernel  $\Theta_{\infty,t}^{(L)}$
- $\Theta_{\infty,t}^{(L)}$  has a **closed-form** expression

# Infinite-Width Limit

**Question:** What happens if we let  $d_1, \dots, d_{L-1} \rightarrow \infty$ ?

We recover the following properties:

- $\Theta_t^{(L)}$  converges to a **deterministic** kernel  $\Theta_{\infty,t}^{(L)}$
- $\Theta_{\infty,t}^{(L)}$  has a **closed-form** expression
- $\Theta_{\infty,t}^{(L)}$  does **not** depend on time:  $\Theta_{\infty,t}^{(L)} = \Theta_{\infty}^{(L)}$

# Infinite-Width Limit

**Question:** What happens if we let  $d_1, \dots, d_{L-1} \rightarrow \infty$ ?

We recover the following properties:

- $\Theta_t^{(L)}$  converges to a **deterministic** kernel  $\Theta_{\infty,t}^{(L)}$
- $\Theta_{\infty,t}^{(L)}$  has a **closed-form** expression
- $\Theta_{\infty,t}^{(L)}$  does **not** depend on time:  $\Theta_{\infty,t}^{(L)} = \Theta_{\infty}^{(L)}$

We hence don't have to worry about the initialization and time-dependence anymore!

# **Consequences**

# Consequences

- Let us **assume** for now that all of this holds. What are the consequences?

# Consequences

- Let us **assume** for now that all of this holds. What are the consequences?
- We have the following dynamics:

## Consequences

- Let us **assume** for now that all of this holds. What are the consequences?
- We have the following dynamics:

$$\partial_t f_t^\infty(\mathbf{x}) = - \left( \boldsymbol{\Theta}_\infty^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T \frac{\partial L(f_t^\infty(\mathbf{X}), \mathbf{y})}{\partial f_t^\infty(\mathbf{X})}$$

# Consequences

- Let us **assume** for now that all of this holds. What are the consequences?
- We have the following dynamics:

$$\partial_t f_t^\infty(\mathbf{x}) = - \left( \boldsymbol{\Theta}_\infty^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T \frac{\partial L(f_t^\infty(\mathbf{X}), \mathbf{y})}{\partial f_t^\infty(\mathbf{X})}$$

- Also assume we use the **squared loss**:

# Consequences

- Let us **assume** for now that all of this holds. What are the consequences?
- We have the following dynamics:

$$\partial_t f_t^\infty(\mathbf{x}) = - \left( \boldsymbol{\Theta}_\infty^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T \frac{\partial L(f_t^\infty(\mathbf{X}), \mathbf{y})}{\partial f_t^\infty(\mathbf{X})}$$

- Also assume we use the **squared loss**:

$$L(f_t^\infty(\mathbf{X}), \mathbf{y}) = \frac{1}{2} \sum_{i=1}^n (f_t^\infty(\mathbf{x}_i) - y_i)^2$$

# Consequences

- Let us **assume** for now that all of this holds. What are the consequences?
- We have the following dynamics:

$$\partial_t f_t^\infty(\mathbf{x}) = - \left( \boldsymbol{\Theta}_\infty^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T \frac{\partial L(f_t^\infty(\mathbf{X}), \mathbf{y})}{\partial f_t^\infty(\mathbf{X})}$$

- Also assume we use the **squared loss**:

$$L(f_t^\infty(\mathbf{X}), \mathbf{y}) = \frac{1}{2} \sum_{i=1}^n (f_t^\infty(\mathbf{x}_i) - y_i)^2$$

- This gives the following expression for the gradient

# Consequences

- Let us **assume** for now that all of this holds. What are the consequences?
- We have the following dynamics:

$$\partial_t f_t^\infty(\mathbf{x}) = - \left( \boldsymbol{\Theta}_\infty^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T \frac{\partial L(f_t^\infty(\mathbf{X}), \mathbf{y})}{\partial f_t^\infty(\mathbf{X})}$$

- Also assume we use the **squared loss**:

$$L(f_t^\infty(\mathbf{X}), \mathbf{y}) = \frac{1}{2} \sum_{i=1}^n (f_t^\infty(\mathbf{x}_i) - y_i)^2$$

- This gives the following expression for the gradient

$$\partial_{f_t^\infty(\mathbf{x}_i)} L(f_t^\infty(\mathbf{X}), \mathbf{y}) = f_t^\infty(\mathbf{x}_i) - y_i$$

# Closed Form Solution (I)

## Closed Form Solution (I)

- Plugging this gradient into our ODE gives:

## Closed Form Solution (I)

- Plugging this gradient into our ODE gives:

$$\partial_t f_t^\infty(\mathbf{x}) = - \left( \boldsymbol{\Theta}_\infty^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T (f_t^\infty(\mathbf{X}) - \mathbf{y})$$

## Closed Form Solution (I)

- Plugging this gradient into our ODE gives:

$$\partial_t f_t^\infty(\mathbf{x}) = - \left( \Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T (f_t^\infty(\mathbf{X}) - \mathbf{y})$$

That's a system of **linear first order** differential equations!

## Closed Form Solution (I)

- Plugging this gradient into our ODE gives:

$$\partial_t f_t^\infty(\mathbf{x}) = - \left( \Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T (f_t^\infty(\mathbf{X}) - \mathbf{y})$$

That's a system of **linear first order** differential equations!

- Careful, every  $\mathbf{x}$  induces a **different** function  $f_t^\infty(\mathbf{x})$  in  $t$ !

## Closed Form Solution (I)

- Plugging this gradient into our ODE gives:

$$\partial_t f_t^\infty(\mathbf{x}) = - \left( \Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T (f_t^\infty(\mathbf{X}) - \mathbf{y})$$

That's a system of **linear first order** differential equations!

- Careful, every  $\mathbf{x}$  induces a **different** function  $f_t^\infty(\mathbf{x})$  in  $t$ !
- After some **tricks** we find a closed form solution:

# Closed Form Solution (I)

- Plugging this gradient into our ODE gives:

$$\partial_t f_t^\infty(\mathbf{x}) = - \left( \Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T (f_t^\infty(\mathbf{X}) - \mathbf{y})$$

That's a system of **linear first order** differential equations!

- Careful, every  $\mathbf{x}$  induces a **different** function  $f_t^\infty(\mathbf{x})$  in  $t$ !
- After some **tricks** we find a closed form solution:

$$f_t^\infty(\mathbf{x}) = \left( \Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T \left( \Theta_\infty^{(L)}(\mathbf{X}, \mathbf{X})^{-1} (e^{-\Theta_\infty^{(L)}(\mathbf{X}, \mathbf{X})t} - 1) \right) (f_0^\infty(\mathbf{X}) - \mathbf{y}) + f_0^\infty(\mathbf{x})$$

# Training an Infinitely-Wide Network Infinitely Long

# Training an Infinitely-Wide Network Infinitely Long

- We can "read off" the **final state** of the infinitely-wide network:

# Training an Infinitely-Wide Network Infinitely Long

- We can "read off" the **final state** of the infinitely-wide network:

$$f_\infty^\infty(\mathbf{x}) = \left( \Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T \left( \Theta_\infty^{(L)}(\mathbf{X}, \mathbf{X}) \right)^{-1} (\mathbf{y} - f_0^\infty(\mathbf{X})) + f_0^\infty(\mathbf{x})$$

# Training an Infinitely-Wide Network Infinitely Long

- We can "read off" the **final state** of the infinitely-wide network:

$$f_\infty^\infty(\mathbf{x}) = \left( \Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T \left( \Theta_\infty^{(L)}(\mathbf{X}, \mathbf{X}) \right)^{-1} (\mathbf{y} - f_0^\infty(\mathbf{X})) + f_0^\infty(\mathbf{x})$$

- We have a **perfect fit** on the training data:

# Training an Infinitely-Wide Network Infinitely Long

- We can "read off" the **final state** of the infinitely-wide network:

$$f_\infty^\infty(\mathbf{x}) = \left( \Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T \left( \Theta_\infty^{(L)}(\mathbf{X}, \mathbf{X}) \right)^{-1} (\mathbf{y} - f_0^\infty(\mathbf{X})) + f_0^\infty(\mathbf{x})$$

- We have a **perfect fit** on the training data:

$$f_\infty^\infty(\mathbf{X}) = \mathbf{y}$$

# Training an Infinitely-Wide Network Infinitely Long

- We can "read off" the **final state** of the infinitely-wide network:

$$f_\infty^\infty(\mathbf{x}) = \left( \Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T \left( \Theta_\infty^{(L)}(\mathbf{X}, \mathbf{X}) \right)^{-1} (\mathbf{y} - f_0^\infty(\mathbf{X})) + f_0^\infty(\mathbf{x})$$

- We have a **perfect fit** on the training data:

$$f_\infty^\infty(\mathbf{X}) = \mathbf{y}$$

(But one has to show that  $\Theta_\infty^{(L)}$  is positive-definite...)

# Training an Infinitely-Wide Network Infinitely Long

- We can "read off" the **final state** of the infinitely-wide network:

$$f_\infty^\infty(\mathbf{x}) = \left( \Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T \left( \Theta_\infty^{(L)}(\mathbf{X}, \mathbf{X}) \right)^{-1} (\mathbf{y} - f_0^\infty(\mathbf{X})) + f_0^\infty(\mathbf{x})$$

- We have a **perfect fit** on the training data:

$$f_\infty^\infty(\mathbf{X}) = \mathbf{y}$$

(But one has to show that  $\Theta_\infty^{(L)}$  is positive-definite...)

- Gradient flow hence converges to the **global minimum** in the infinite regime.

# Kernel Ridge Regression

# Kernel Ridge Regression

- If we assume  $f_0^\infty(\mathbf{x}) = 0$ , we have

# Kernel Ridge Regression

- If we assume  $f_0^\infty(\mathbf{x}) = 0$ , we have

$$f_\infty^\infty(\mathbf{x}) = \mathbf{y}^T \left( \Theta_\infty^{(L)}(\mathbf{X}, \mathbf{X}) \right)^{-1} \Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X})$$

# Kernel Ridge Regression

- If we assume  $f_0^\infty(\mathbf{x}) = 0$ , we have

$$f_\infty^\infty(\mathbf{x}) = \mathbf{y}^T \left( \Theta_\infty^{(L)}(\mathbf{X}, \mathbf{X}) \right)^{-1} \Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X})$$

- That's exactly the same solution as in unregularized **kernel regression** with a positive-definite kernel K:

# Kernel Ridge Regression

- If we assume  $f_0^\infty(\mathbf{x}) = 0$ , we have

$$f_\infty^\infty(\mathbf{x}) = \mathbf{y}^T \left( \Theta_\infty^{(L)}(\mathbf{X}, \mathbf{X}) \right)^{-1} \Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X})$$

- That's exactly the same solution as in unregularized **kernel regression** with a positive-definite kernel K:

$$\hat{y} = \mathbf{y}^T K(\mathbf{X}, \mathbf{X})^{-1} K(\mathbf{x}, \mathbf{X})$$

# Kernel Ridge Regression

- If we assume  $f_0^\infty(\mathbf{x}) = 0$ , we have

$$f_\infty^\infty(\mathbf{x}) = \mathbf{y}^T \left( \Theta_\infty^{(L)}(\mathbf{X}, \mathbf{X}) \right)^{-1} \Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X})$$

- That's exactly the same solution as in unregularized **kernel regression** with a positive-definite kernel K:

$$\hat{\mathbf{y}} = \mathbf{y}^T K(\mathbf{X}, \mathbf{X})^{-1} K(\mathbf{x}, \mathbf{X})$$

- Learning an infinitely-wide neural network **is the same** as kernel regression with  $\Theta_\infty^{(L)}$

# **Let's Reiterate**

## Let's Reiterate

- $\Theta_t^{(L)}(\mathbf{x}, \mathbf{x}') = \sum_{p=1}^P \partial_{\theta_p} f_{\theta}(\mathbf{x}) \partial_{\theta_p} f_{\theta}(\mathbf{x}')$

## Let's Reiterate

- $\Theta_t^{(L)}(\mathbf{x}, \mathbf{x}') = \sum_{p=1}^P \partial_{\theta_p} f_{\theta}(\mathbf{x}) \partial_{\theta_p} f_{\theta}(\mathbf{x}')$
- $\partial_t f_{\theta}(\mathbf{x}) = -\Theta_t^{(L)}(\mathbf{x}, \mathbf{X})^T \frac{\partial L(f_{\theta}(\mathbf{X}), \mathbf{y})}{\partial f_{\theta}(\mathbf{X})}$

## Let's Reiterate

- $\Theta_t^{(L)}(\mathbf{x}, \mathbf{x}') = \sum_{p=1}^P \partial_{\theta_p} f_{\theta}(\mathbf{x}) \partial_{\theta_p} f_{\theta}(\mathbf{x}')$
- $\partial_t f_{\theta}(\mathbf{x}) = -\Theta_t^{(L)}(\mathbf{x}, \mathbf{X})^T \frac{\partial L(f_{\theta}(\mathbf{X}), \mathbf{y})}{\partial f_{\theta}(\mathbf{X})}$
- $\Theta_{\infty}^{(L)}$  closed-form, deterministic and time-independent

## Let's Reiterate

- $\Theta_t^{(L)}(\mathbf{x}, \mathbf{x}') = \sum_{p=1}^P \partial_{\theta_p} f_{\theta}(\mathbf{x}) \partial_{\theta_p} f_{\theta}(\mathbf{x}')$
  - $\partial_t f_{\theta}(\mathbf{x}) = -\Theta_t^{(L)}(\mathbf{x}, \mathbf{X})^T \frac{\partial L(f_{\theta}(\mathbf{X}), \mathbf{y})}{\partial f_{\theta}(\mathbf{X})}$
  - $\Theta_{\infty}^{(L)}$  closed-form, deterministic and time-independent
  - Recover final state of the infinite model
- $$f_{\infty}^{\infty}(\mathbf{x}) = \left( \Theta_{\infty}^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T \left( \Theta_{\infty}^{(L)}(\mathbf{X}, \mathbf{X}) \right)^{-1} \left( \mathbf{y} - f_0^{\infty}(\mathbf{X}) \right) + f_0^{\infty}(\mathbf{x})$$

## Let's Reiterate

- $\Theta_t^{(L)}(\mathbf{x}, \mathbf{x}') = \sum_{p=1}^P \partial_{\theta_p} f_{\theta}(\mathbf{x}) \partial_{\theta_p} f_{\theta}(\mathbf{x}')$
  - $\partial_t f_{\theta}(\mathbf{x}) = -\Theta_t^{(L)}(\mathbf{x}, \mathbf{X})^T \frac{\partial L(f_{\theta}(\mathbf{X}), \mathbf{y})}{\partial f_{\theta}(\mathbf{X})}$
  - $\Theta_{\infty}^{(L)}$  closed-form, deterministic and time-independent
  - Recover final state of the infinite model
- $$f_{\infty}^{\infty}(\mathbf{x}) = \left( \Theta_{\infty}^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T \left( \Theta_{\infty}^{(L)}(\mathbf{X}, \mathbf{X}) \right)^{-1} \left( \mathbf{y} - f_0^{\infty}(\mathbf{X}) \right) + f_0^{\infty}(\mathbf{x})$$

# Recap: NNs as GPs at Initialization

**Theorem 2:**

## Recap: NNs as GPs at Initialization

### Theorem 2:

Given a NN  $f_\theta$ , it holds that for  $d_1, \dots, d_{L-1} \rightarrow \infty$ :

$$f_{\theta,i}^\infty(\cdot) \sim \mathcal{GP}(0, \Sigma^{(L)})$$

# Recap: NNs as GPs at Initialization

## Theorem 2:

Given a NN  $f_\theta$ , it holds that for  $d_1, \dots, d_{L-1} \rightarrow \infty$ :

$$f_{\theta,i}^\infty(\cdot) \sim \mathcal{GP}(0, \Sigma^{(L)})$$

i.i.d. in  $i$ , where  $\Sigma^{(L)}$  has the recursive structure

- $\Sigma^{(1)}(\mathbf{x}, \mathbf{x}') = \frac{1}{\sqrt{d_0}} \mathbf{x}^T \mathbf{x}' + \beta^2$
- $\Sigma^{(l)}(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \tilde{\Sigma}^{(l-1)})} [\sigma(z_1)\sigma(z_2)] + \beta^2$

# Recap: NNs as GPs at Initialization

## Theorem 2:

Given a NN  $f_\theta$ , it holds that for  $d_1, \dots, d_{L-1} \rightarrow \infty$ :

$$f_{\theta,i}^\infty(\cdot) \sim \mathcal{GP}(0, \Sigma^{(L)})$$

i.i.d. in  $i$ , where  $\Sigma^{(L)}$  has the recursive structure

- $\Sigma^{(1)}(\mathbf{x}, \mathbf{x}') = \frac{1}{\sqrt{d_0}} \mathbf{x}^T \mathbf{x}' + \beta^2$
- $\Sigma^{(l)}(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \tilde{\Sigma}^{(l-1)})} [\sigma(z_1)\sigma(z_2)] + \beta^2$

$$\text{for } \tilde{\Sigma}^{(l-1)} = \begin{pmatrix} \Sigma^{(l-1)}(\mathbf{x}, \mathbf{x}) & \Sigma^{(l-1)}(\mathbf{x}, \mathbf{x}') \\ \Sigma^{(l-1)}(\mathbf{x}', \mathbf{x}) & \Sigma^{(l-1)}(\mathbf{x}', \mathbf{x}') \end{pmatrix}$$

# NTK at Initialization

**Theorem 3:**

# NTK at Initialization

## Theorem 3:

*Take a NN  $f_\theta$  of depth  $L$  with Lipschitz  $\sigma$ . Then, for  $d_1, \dots, d_L \rightarrow \infty$ , it holds*

# NTK at Initialization

## Theorem 3:

Take a NN  $f_\theta$  of depth  $L$  with Lipschitz  $\sigma$ . Then, for  $d_1, \dots, d_L \rightarrow \infty$ , it holds

$$\Theta_0^{(L)} \xrightarrow{\mathbb{P}} \Theta_{0,\infty}^{(L)} \otimes \mathbb{1}_{d_L \times d_L}$$

# NTK at Initialization

## Theorem 3:

Take a NN  $f_\theta$  of depth  $L$  with Lipschitz  $\sigma$ . Then, for  $d_1, \dots, d_L \rightarrow \infty$ , it holds

$$\Theta_0^{(L)} \xrightarrow{\mathbb{P}} \Theta_{0,\infty}^{(L)} \otimes \mathbb{1}_{d_L \times d_L}$$

where  $\Theta_{0,\infty}^{(L)} : \mathbb{R}^{d_0} \times \mathbb{R}^{d_0} \rightarrow \mathbb{R}$  follows the recursion

# NTK at Initialization

## Theorem 3:

Take a NN  $f_\theta$  of depth  $L$  with Lipschitz  $\sigma$ . Then, for  $d_1, \dots, d_L \rightarrow \infty$ , it holds

$$\Theta_0^{(L)} \xrightarrow{\mathbb{P}} \Theta_{0,\infty}^{(L)} \otimes \mathbb{1}_{d_L \times d_L}$$

where  $\Theta_{0,\infty}^{(L)} : \mathbb{R}^{d_0} \times \mathbb{R}^{d_0} \rightarrow \mathbb{R}$  follows the recursion

- $\Theta_{0,\infty}^{(1)}(\mathbf{x}, \mathbf{x}') = \Sigma^{(1)}(\mathbf{x}, \mathbf{x}')$
- $\Theta_{0,\infty}^{(l+1)}(\mathbf{x}, \mathbf{x}') = \Theta_{0,\infty}^{(l)}(\mathbf{x}, \mathbf{x}') \dot{\Sigma}^{(l+1)}(\mathbf{x}, \mathbf{x}') + \Sigma^{(l+1)}(\mathbf{x}, \mathbf{x}')$

# NTK at Initialization

## Theorem 3:

Take a NN  $f_\theta$  of depth  $L$  with Lipschitz  $\sigma$ . Then, for  $d_1, \dots, d_L \rightarrow \infty$ , it holds

$$\Theta_0^{(L)} \xrightarrow{\mathbb{P}} \Theta_{0,\infty}^{(L)} \otimes \mathbb{1}_{d_L \times d_L}$$

where  $\Theta_{0,\infty}^{(L)} : \mathbb{R}^{d_0} \times \mathbb{R}^{d_0} \rightarrow \mathbb{R}$  follows the recursion

- $\Theta_{0,\infty}^{(1)}(\mathbf{x}, \mathbf{x}') = \Sigma^{(1)}(\mathbf{x}, \mathbf{x}')$
- $\Theta_{0,\infty}^{(l+1)}(\mathbf{x}, \mathbf{x}') = \Theta_{0,\infty}^{(l)}(\mathbf{x}, \mathbf{x}') \dot{\Sigma}^{(l+1)}(\mathbf{x}, \mathbf{x}') + \Sigma^{(l+1)}(\mathbf{x}, \mathbf{x}')$

$$\text{with } \dot{\Sigma}^{(l)}(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \tilde{\Sigma}^{(l)})} \left[ \dot{\sigma}(z_1) \dot{\sigma}(z_2) \right]$$

# NTK at Initialization

## Theorem 3:

Take a NN  $f_\theta$  of depth  $L$  with Lipschitz  $\sigma$ . Then, for  $d_1, \dots, d_L \rightarrow \infty$ , it holds

$$\Theta_0^{(L)} \xrightarrow{\mathbb{P}} \Theta_{0,\infty}^{(L)} \otimes \mathbb{1}_{d_L \times d_L}$$

where  $\Theta_{0,\infty}^{(L)} : \mathbb{R}^{d_0} \times \mathbb{R}^{d_0} \rightarrow \mathbb{R}$  follows the recursion

- $\Theta_{0,\infty}^{(1)}(\mathbf{x}, \mathbf{x}') = \Sigma^{(1)}(\mathbf{x}, \mathbf{x}')$
- $\Theta_{0,\infty}^{(l+1)}(\mathbf{x}, \mathbf{x}') = \Theta_{0,\infty}^{(l)}(\mathbf{x}, \mathbf{x}') \dot{\Sigma}^{(l+1)}(\mathbf{x}, \mathbf{x}') + \Sigma^{(l+1)}(\mathbf{x}, \mathbf{x}')$

$$\text{with } \dot{\Sigma}^{(l)}(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \tilde{\Sigma}^{(l)})} \left[ \dot{\sigma}(z_1) \dot{\sigma}(z_2) \right]$$

Writing the expression on the scalar level gives

# NTK at Initialization

## Theorem 3:

Take a NN  $f_\theta$  of depth  $L$  with Lipschitz  $\sigma$ . Then, for  $d_1, \dots, d_L \rightarrow \infty$ , it holds

$$\Theta_0^{(L)} \xrightarrow{\mathbb{P}} \Theta_{0,\infty}^{(L)} \otimes \mathbb{1}_{d_L \times d_L}$$

where  $\Theta_{0,\infty}^{(L)} : \mathbb{R}^{d_0} \times \mathbb{R}^{d_0} \rightarrow \mathbb{R}$  follows the recursion

- $\Theta_{0,\infty}^{(1)}(\mathbf{x}, \mathbf{x}') = \Sigma^{(1)}(\mathbf{x}, \mathbf{x}')$
- $\Theta_{0,\infty}^{(l+1)}(\mathbf{x}, \mathbf{x}') = \Theta_{0,\infty}^{(l)}(\mathbf{x}, \mathbf{x}') \dot{\Sigma}^{(l+1)}(\mathbf{x}, \mathbf{x}') + \Sigma^{(l+1)}(\mathbf{x}, \mathbf{x}')$

$$\text{with } \dot{\Sigma}^{(l)}(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \tilde{\Sigma}^{(l)})} \left[ \dot{\sigma}(z_1) \dot{\sigma}(z_2) \right]$$

Writing the expression on the scalar level gives

$$\left( \Theta_{0,\infty}^{(L)} \otimes \mathbb{1}_{d_L \times d_L} \right)_{kk'}(\mathbf{x}, \mathbf{x}') = \Theta_{0,\infty}^{(L)}(\mathbf{x}, \mathbf{x}') \mathbb{1}_{\{k=k'\}}$$

# **Proof**

# Proof

- Proof is based on **induction** on the depth  $L$

# Proof

- Proof is based on **induction** on the depth  $L$
- Limits  $d_1, \dots, d_{L-1} \rightarrow \infty$  again **sequentially** taken

# Proof

- Proof is based on **induction** on the depth  $L$
- Limits  $d_1, \dots, d_{L-1} \rightarrow \infty$  again **sequentially** taken
- For ease of notation, assume  $\beta = 0$ , so **no bias**

# Proof

- Proof is based on **induction** on the depth  $L$
- Limits  $d_1, \dots, d_{L-1} \rightarrow \infty$  again **sequentially** taken
- For ease of notation, assume  $\beta = 0$ , so **no bias**
- Can't do it for just  $d_L = 1$ , otherwise induction doesn't work.  
Have to do it for **general**  $d_L\dots$

**Base Case:**  $L = 1$  (I)

## **Base Case: $L = 1$ (I)**

In this case the network is given by

## Base Case: $L = 1$ (I)

In this case the network is given by

$$f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_0}} \sum_{l=1}^{d_0} W_{kl}^{(0)} x_l$$

## Base Case: $L = 1$ (I)

In this case the network is given by

$$f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_0}} \sum_{l=1}^{d_0} W_{kl}^{(0)} x_l$$

Gradients can hence be written as

## Base Case: $L = 1$ (I)

In this case the network is given by

$$f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_0}} \sum_{l=1}^{d_0} W_{kl}^{(0)} x_l$$

Gradients can hence be written as

$$\partial_{W_{ij}^{(0)}} f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_0}} x_j \mathbb{1}_{\{k=i\}}$$

## Base Case: $L = 1$ (I)

In this case the network is given by

$$f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_0}} \sum_{l=1}^{d_0} W_{kl}^{(0)} x_l$$

Gradients can hence be written as

$$\partial_{W_{ij}^{(0)}} f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_0}} x_j \mathbb{1}_{\{k=i\}}$$

**Remember:**  $\Theta_{kk'}^{(L)}(\mathbf{x}, \mathbf{x}') = \sum_{p=1}^P \partial_{\theta_p} f_{\theta,k}(\mathbf{x}) \partial_{\theta_p} f_{\theta,k'}(\mathbf{x}')$

**Base Case:**  $L = 1$  (II)

## **Base Case: $L = 1$ (II)**

Plugging in the definition of the NTK leads to:

## Base Case: $L = 1$ (II)

Plugging in the definition of the NTK leads to:

$$\Theta_{k,k'}^{(1)}(\mathbf{x}, \mathbf{x}')$$

## Base Case: $L = 1$ (II)

Plugging in the definition of the NTK leads to:

$$\Theta_{k,k'}^{(1)}(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^{d_0} \sum_{j=1}^{d_1} \partial_{W_{ij}^{(0)}} f_{\theta,k}(\mathbf{x}) \partial_{W_{ij}^{(0)}} f_{\theta,k'}(\mathbf{x}')$$

## Base Case: $L = 1$ (II)

Plugging in the definition of the NTK leads to:

$$\begin{aligned}\Theta_{k,k'}^{(1)}(\mathbf{x}, \mathbf{x}') &= \sum_{i=1}^{d_0} \sum_{j=1}^{d_1} \partial_{W_{ij}^{(0)}} f_{\theta,k}(\mathbf{x}) \partial_{W_{ij}^{(0)}} f_{\theta,k'}(\mathbf{x}') \\ &= \sum_{i=1}^{d_0} \sum_{j=1}^{d_1} \frac{1}{d_0} x_j \mathbb{1}_{\{k=i\}} x'_j \mathbb{1}_{\{k'=i\}}\end{aligned}$$

## Base Case: $L = 1$ (II)

Plugging in the definition of the NTK leads to:

$$\begin{aligned}\Theta_{k,k'}^{(1)}(\mathbf{x}, \mathbf{x}') &= \sum_{i=1}^{d_0} \sum_{j=1}^{d_1} \partial_{W_{ij}^{(0)}} f_{\theta,k}(\mathbf{x}) \partial_{W_{ij}^{(0)}} f_{\theta,k'}(\mathbf{x}') \\ &= \sum_{i=1}^{d_0} \sum_{j=1}^{d_1} \frac{1}{d_0} x_j \mathbb{1}_{\{k=i\}} x'_j \mathbb{1}_{\{k'=i\}} \\ &= \frac{1}{d_0} \mathbf{x}^T \mathbf{x}' \mathbb{1}_{\{k'=k\}}\end{aligned}$$

## Base Case: $L = 1$ (II)

Plugging in the definition of the NTK leads to:

$$\begin{aligned}\Theta_{k,k'}^{(1)}(\mathbf{x}, \mathbf{x}') &= \sum_{i=1}^{d_0} \sum_{j=1}^{d_1} \partial_{W_{ij}^{(0)}} f_{\theta,k}(\mathbf{x}) \partial_{W_{ij}^{(0)}} f_{\theta,k'}(\mathbf{x}') \\ &= \sum_{i=1}^{d_0} \sum_{j=1}^{d_1} \frac{1}{d_0} x_j \mathbb{1}_{\{k=i\}} x'_j \mathbb{1}_{\{k'=i\}} \\ &= \frac{1}{d_0} \mathbf{x}^T \mathbf{x}' \mathbb{1}_{\{k'=k\}} \\ &= \Sigma^{(1)}(\mathbf{x}, \mathbf{x}') \mathbb{1}_{\{k'=k\}}\end{aligned}$$

**Notice:** We didn't take any limit because input dimension  $d_0$  of course remains fixed.

## Base Case: $L = 1$ (II)

Plugging in the definition of the NTK leads to:

$$\begin{aligned}\Theta_{k,k'}^{(1)}(\mathbf{x}, \mathbf{x}') &= \sum_{i=1}^{d_0} \sum_{j=1}^{d_1} \partial_{W_{ij}^{(0)}} f_{\theta,k}(\mathbf{x}) \partial_{W_{ij}^{(0)}} f_{\theta,k'}(\mathbf{x}') \\ &= \sum_{i=1}^{d_0} \sum_{j=1}^{d_1} \frac{1}{d_0} x_j \mathbb{1}_{\{k=i\}} x'_j \mathbb{1}_{\{k'=i\}} \\ &= \frac{1}{d_0} \mathbf{x}^T \mathbf{x}' \mathbb{1}_{\{k'=k\}} \\ &= \Sigma^{(1)}(\mathbf{x}, \mathbf{x}') \mathbb{1}_{\{k'=k\}}\end{aligned}$$

**Notice:** We didn't take any limit because input dimension  $d_0$  of course remains fixed.

This **proves** the base case.

# **Induction Step**

## Induction Step

- Take a network of depth  $l + 1$ :

## Induction Step

- Take a network of depth  $l + 1$ :

$$f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_l}} \sum_{j=1}^{d_l} W_{kj}^{(l)} \sigma \left( \tilde{\alpha}_j^{(l)}(\mathbf{x}) \right)$$

## Induction Step

- Take a network of depth  $l + 1$ :

$$f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_l}} \sum_{j=1}^{d_l} W_{kj}^{(l)} \sigma \left( \tilde{\alpha}_j^{(l)}(\mathbf{x}) \right)$$

- Split the parameters into  $\theta = (\tilde{\theta}, \mathbf{W}^{(l)})$

# Induction Step

- Take a network of depth  $l + 1$ :

$$f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_l}} \sum_{j=1}^{d_l} W_{kj}^{(l)} \sigma \left( \tilde{\alpha}_j^{(l)}(\mathbf{x}) \right)$$

- Split the parameters into  $\theta = (\tilde{\theta}, \mathbf{W}^{(l)})$
- We know that  $\Theta_{kk'}^{(l)}(\mathbf{x}, \mathbf{x}') \xrightarrow{d_1, \dots, d_{l-1} \rightarrow \infty} \Theta_\infty^{(l)}(\mathbf{x}, \mathbf{x}') \mathbb{1}_{\{k=k'\}}$

## Induction Step

- Take a network of depth  $l + 1$ :

$$f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_l}} \sum_{j=1}^{d_l} W_{kj}^{(l)} \sigma \left( \tilde{\alpha}_j^{(l)}(\mathbf{x}) \right)$$

- Split the parameters into  $\theta = (\tilde{\theta}, \mathbf{W}^{(l)})$
- We know that  $\Theta_{kk'}^{(l)}(\mathbf{x}, \mathbf{x}') \xrightarrow{d_1, \dots, d_{l-1} \rightarrow \infty} \Theta_\infty^{(l)}(\mathbf{x}, \mathbf{x}') \mathbb{1}_{\{k=k'\}}$

We can now **split** the NTK:

# Induction Step

- Take a network of depth  $l + 1$ :

$$f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_l}} \sum_{j=1}^{d_l} W_{kj}^{(l)} \sigma \left( \tilde{\alpha}_j^{(l)}(\mathbf{x}) \right)$$

- Split the parameters into  $\theta = (\tilde{\theta}, \mathbf{W}^{(l)})$
- We know that  $\Theta_{kk'}^{(l)}(\mathbf{x}, \mathbf{x}') \xrightarrow{d_1, \dots, d_{l-1} \rightarrow \infty} \Theta_\infty^{(l)}(\mathbf{x}, \mathbf{x}') \mathbb{1}_{\{k=k'\}}$

We can now **split** the NTK:

$$\Theta_{kk'}^{(l+1)}(\mathbf{x}, \mathbf{x}') = \underbrace{\sum_{l=1}^{\tilde{P}} \partial_{\tilde{\theta}_l} f_{\theta,k}(\mathbf{x}), \partial_{\tilde{\theta}_l} f_{\theta,k'}(\mathbf{x}')}_{A} + \underbrace{\sum_{i,j=1}^{d_l, d_{l+1}} \partial_{W_{ij}^{(l)}} f_{\theta,k}(\mathbf{x}), \partial_{W_{ij}^{(l)}} f_{\theta,k'}(\mathbf{x}')}_{B}$$

# Terms in A (I)

## Terms in A (I)

- Write out the derivatives as

## Terms in A (I)

- Write out the derivatives as

$$\partial_{\tilde{\theta}_p} f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_L}} \sum_{j=1}^{d_L} W_{kj}^{(I)} \dot{\sigma}\left(\tilde{\alpha}_j(\mathbf{x})\right) \partial_{\tilde{\theta}_p} \tilde{\alpha}_j(\mathbf{x})$$

## Terms in A (I)

- Write out the derivatives as

$$\partial_{\tilde{\theta}_p} f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_L}} \sum_{j=1}^{d_L} W_{kj}^{(l)} \dot{\sigma}\left(\tilde{\alpha}_j(\mathbf{x})\right) \partial_{\tilde{\theta}_p} \tilde{\alpha}_j(\mathbf{x})$$

Let us study the terms  $A = \sum_{l=1}^{\tilde{P}} \partial_{\tilde{\theta}_l} f_{\theta,k}(\mathbf{x}) \partial_{\tilde{\theta}_l} f_{\theta,k'}(\mathbf{x}')$  due to  $\tilde{\theta}$ :

## Terms in A (I)

- Write out the derivatives as

$$\partial_{\tilde{\theta}_p} f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_L}} \sum_{j=1}^{d_L} W_{kj}^{(l)} \dot{\sigma}\left(\tilde{\alpha}_j(\mathbf{x})\right) \partial_{\tilde{\theta}_p} \tilde{\alpha}_j(\mathbf{x})$$

Let us study the terms  $A = \sum_{l=1}^{\tilde{P}} \partial_{\tilde{\theta}_l} f_{\theta,k}(\mathbf{x}) \partial_{\tilde{\theta}_l} f_{\theta,k'}(\mathbf{x}')$  due to  $\tilde{\theta}$ :

$$\frac{1}{d_l} \sum_{p=1}^{\tilde{P}} \left( \sum_{j=1}^{d_l} W_{kj}^{(l)} \dot{\sigma}\left(\tilde{\alpha}_j^{(l)}(\mathbf{x})\right) \partial_{\tilde{\theta}_p} \tilde{\alpha}_j^{(l)}(\mathbf{x}) \right) \left( \sum_{j'=1}^{d_l} W_{k'j'}^{(l)} \dot{\sigma}\left(\tilde{\alpha}_{j'}^{(l)}(\mathbf{x}')\right) \partial_{\tilde{\theta}_p} \alpha_{j'}^{(l)}(\mathbf{x}') \right)$$

## Terms in A (I)

- Write out the derivatives as

$$\partial_{\tilde{\theta}_p} f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_L}} \sum_{j=1}^{d_L} W_{kj}^{(l)} \dot{\sigma}\left(\tilde{\alpha}_j(\mathbf{x})\right) \partial_{\tilde{\theta}_p} \tilde{\alpha}_j(\mathbf{x})$$

Let us study the terms  $A = \sum_{l=1}^{\tilde{P}} \partial_{\tilde{\theta}_l} f_{\theta,k}(\mathbf{x}) \partial_{\tilde{\theta}_l} f_{\theta,k'}(\mathbf{x}')$  due to  $\tilde{\theta}$ :

$$\begin{aligned} & \frac{1}{d_l} \sum_{p=1}^{\tilde{P}} \left( \sum_{j=1}^{d_l} W_{kj}^{(l)} \dot{\sigma}\left(\tilde{\alpha}_j^{(l)}(\mathbf{x})\right) \partial_{\tilde{\theta}_p} \tilde{\alpha}_j^{(l)}(\mathbf{x}) \right) \left( \sum_{j'=1}^{d_l} W_{k'j'}^{(l)} \dot{\sigma}\left(\tilde{\alpha}_{j'}^{(l)}(\mathbf{x}')\right) \partial_{\tilde{\theta}_p} \alpha_{j'}^{(l)}(\mathbf{x}') \right) \\ &= \frac{1}{d_l} \sum_{j,j'=1}^{d_l} W_{kj}^{(l)} W_{k'j'}^{(l)} \dot{\sigma}\left(\tilde{\alpha}_j^{(l)}(\mathbf{x})\right) \dot{\sigma}\left(\tilde{\alpha}_{j'}^{(l)}(\mathbf{x}')\right) \sum_{p=1}^{\tilde{P}} \partial_{\tilde{\theta}_p} \tilde{\alpha}_j^{(l)}(\mathbf{x}) \partial_{\tilde{\theta}_p} \tilde{\alpha}_{j'}^{(l)}(\mathbf{x}') \end{aligned}$$

# Terms in A (I)

- Write out the derivatives as

$$\partial_{\tilde{\theta}_p} f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_L}} \sum_{j=1}^{d_L} W_{kj}^{(l)} \dot{\sigma}\left(\tilde{\alpha}_j(\mathbf{x})\right) \partial_{\tilde{\theta}_p} \tilde{\alpha}_j(\mathbf{x})$$

Let us study the terms  $A = \sum_{l=1}^{\tilde{P}} \partial_{\tilde{\theta}_l} f_{\theta,k}(\mathbf{x}) \partial_{\tilde{\theta}_l} f_{\theta,k'}(\mathbf{x}')$  due to  $\tilde{\theta}$ :

$$\begin{aligned} & \frac{1}{d_l} \sum_{p=1}^{\tilde{P}} \left( \sum_{j=1}^{d_l} W_{kj}^{(l)} \dot{\sigma}\left(\tilde{\alpha}_j^{(l)}(\mathbf{x})\right) \partial_{\tilde{\theta}_p} \tilde{\alpha}_j^{(l)}(\mathbf{x}) \right) \left( \sum_{j'=1}^{d_l} W_{k'j'}^{(l)} \dot{\sigma}\left(\tilde{\alpha}_{j'}^{(l)}(\mathbf{x}')\right) \partial_{\tilde{\theta}_p} \alpha_{j'}^{(l)}(\mathbf{x}') \right) \\ &= \frac{1}{d_l} \sum_{j,j'=1}^{d_l} W_{kj}^{(l)} W_{k'j'}^{(l)} \dot{\sigma}\left(\tilde{\alpha}_j^{(l)}(\mathbf{x})\right) \dot{\sigma}\left(\tilde{\alpha}_{j'}^{(l)}(\mathbf{x}')\right) \sum_{p=1}^{\tilde{P}} \partial_{\tilde{\theta}_p} \tilde{\alpha}_j^{(l)}(\mathbf{x}) \partial_{\tilde{\theta}_p} \tilde{\alpha}_{j'}^{(l)}(\mathbf{x}') \\ &= \frac{1}{d_l} \sum_{j,j'=1}^{d_l} W_{kj}^{(l)} W_{k'j'}^{(l)} \dot{\sigma}\left(\tilde{\alpha}_j^{(l)}(\mathbf{x})\right) \dot{\sigma}\left(\tilde{\alpha}_{j'}^{(l)}(\mathbf{x}')\right) \Theta_{j,j'}^{(l)}(\mathbf{x}, \mathbf{x}') \end{aligned}$$

## **Terms in A (II)**

## Terms in A (II)

$A$

## Terms in A (II)

$$A = \frac{1}{d_l} \sum_{j,j'=1}^{d_l} W_{kj}^{(l)} W_{k'j'}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_j^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_{j'}^{(l)}(\mathbf{x}') \right) \Theta_{j,j'}^{(l)}(\mathbf{x}, \mathbf{x}')$$

## Terms in A (II)

$$A = \frac{1}{d_l} \sum_{j,j'=1}^{d_l} W_{kj}^{(l)} W_{k'j'}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_j^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_{j'}^{(l)}(\mathbf{x}') \right) \Theta_{j,j'}^{(l)}(\mathbf{x}, \mathbf{x}')$$

$$\xrightarrow{d_1, \dots, d_{l-1} \rightarrow \infty} \frac{1}{d_l} \sum_{j,j'=1}^{d_l} W_{kj}^{(l)} W_{k'j'}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_{j',\infty}^{(l)}(\mathbf{x}') \right) \Theta_{\infty}^{(l)}(\mathbf{x}, \mathbf{x}') \mathbb{1}_{\{j=j'\}}$$

## Terms in A (II)

$$A = \frac{1}{d_l} \sum_{j,j'=1}^{d_l} W_{kj}^{(l)} W_{k'j'}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_j^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_{j'}^{(l)}(\mathbf{x}') \right) \Theta_{j,j'}^{(l)}(\mathbf{x}, \mathbf{x}')$$

$$\xrightarrow{d_1, \dots, d_{l-1} \rightarrow \infty} \frac{1}{d_l} \sum_{j,j'=1}^{d_l} W_{kj}^{(l)} W_{k'j'}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_{j',\infty}^{(l)}(\mathbf{x}') \right) \Theta_{\infty}^{(l)}(\mathbf{x}, \mathbf{x}') \mathbb{1}_{\{j=j'\}}$$

$$= \frac{1}{d_l} \sum_{j=1}^{d_l} W_{kj}^{(l)} W_{k'j}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}') \right) \Theta_{\infty}^{(l)}(\mathbf{x}, \mathbf{x}')$$

## Terms in A (II)

$$A = \frac{1}{d_l} \sum_{j,j'=1}^{d_l} W_{kj}^{(l)} W_{k'j'}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_j^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_{j'}^{(l)}(\mathbf{x}') \right) \Theta_{j,j'}^{(l)}(\mathbf{x}, \mathbf{x}')$$

$$\xrightarrow{d_1, \dots, d_{l-1} \rightarrow \infty} \frac{1}{d_l} \sum_{j,j'=1}^{d_l} W_{kj}^{(l)} W_{k'j'}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_{j',\infty}^{(l)}(\mathbf{x}') \right) \Theta_{\infty}^{(l)}(\mathbf{x}, \mathbf{x}') \mathbb{1}_{\{j=j'\}}$$

$$= \frac{1}{d_l} \sum_{j=1}^{d_l} W_{kj}^{(l)} W_{k'j}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}') \right) \Theta_{\infty}^{(l)}(\mathbf{x}, \mathbf{x}')$$

$$\xrightarrow{d_l \rightarrow \infty} \mathbb{E} \left[ W_{kj}^{(L)} W_{k'j}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}') \right) \Theta_{\infty}^{(l)}(\mathbf{x}, \mathbf{x}') \right]$$

## Terms in A (II)

$$A = \frac{1}{d_l} \sum_{j,j'=1}^{d_l} W_{kj}^{(l)} W_{k'j'}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_j^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_{j'}^{(l)}(\mathbf{x}') \right) \Theta_{j,j'}^{(l)}(\mathbf{x}, \mathbf{x}')$$

$$\xrightarrow{d_1, \dots, d_{l-1} \rightarrow \infty} \frac{1}{d_l} \sum_{j,j'=1}^{d_l} W_{kj}^{(l)} W_{k'j'}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_{j',\infty}^{(l)}(\mathbf{x}') \right) \Theta_{\infty}^{(l)}(\mathbf{x}, \mathbf{x}') \mathbb{1}_{\{j=j'\}}$$

$$= \frac{1}{d_l} \sum_{j=1}^{d_l} W_{kj}^{(l)} W_{k'j}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}') \right) \Theta_{\infty}^{(l)}(\mathbf{x}, \mathbf{x}')$$

$$\xrightarrow{d_l \rightarrow \infty} \mathbb{E} \left[ W_{kj}^{(L)} W_{k'j}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}') \right) \Theta_{\infty}^{(l)}(\mathbf{x}, \mathbf{x}') \right]$$

$$= \Theta_{\infty}^{(L)}(\mathbf{x}, \mathbf{x}') \mathbb{E} \left[ W_{kj}^{(l)} W_{k'j}^{(l)} \right] \mathbb{E} \left[ \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}, \tilde{\theta}) \right) \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}') \right) \right]$$

# Terms in A (II)

$$A = \frac{1}{d_l} \sum_{j,j'=1}^{d_l} W_{kj}^{(l)} W_{k'j'}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_j^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_{j'}^{(l)}(\mathbf{x}') \right) \Theta_{j,j'}^{(l)}(\mathbf{x}, \mathbf{x}')$$

$$\xrightarrow{d_1, \dots, d_{l-1} \rightarrow \infty} \frac{1}{d_l} \sum_{j,j'=1}^{d_l} W_{kj}^{(l)} W_{k'j'}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_{j',\infty}^{(l)}(\mathbf{x}') \right) \Theta_{\infty}^{(l)}(\mathbf{x}, \mathbf{x}') \mathbb{1}_{\{j=j'\}}$$

$$= \frac{1}{d_l} \sum_{j=1}^{d_l} W_{kj}^{(l)} W_{k'j}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}') \right) \Theta_{\infty}^{(l)}(\mathbf{x}, \mathbf{x}')$$

$$\xrightarrow{d_l \rightarrow \infty} \mathbb{E} \left[ W_{kj}^{(L)} W_{k'j}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}') \right) \Theta_{\infty}^{(l)}(\mathbf{x}, \mathbf{x}') \right]$$

$$= \Theta_{\infty}^{(L)}(\mathbf{x}, \mathbf{x}') \mathbb{E} \left[ W_{kj}^{(l)} W_{k'j}^{(l)} \right] \mathbb{E} \left[ \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}, \tilde{\theta}) \right) \dot{\sigma} \left( \tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}') \right) \right]$$

$$= \Theta_{\infty}^{(l)}(\mathbf{x}, \mathbf{x}') \mathbb{1}_{\{k=k'\}} \dot{\Sigma}^{(L+1)}(\mathbf{x}, \mathbf{x}')$$

## **Terms in B**

## Terms in B

Let's calculate the gradients w.r.t. last layer weights

# Terms in B

Let's calculate the gradients w.r.t. last layer weights

$$\partial_{W_{ij}^{(l)}} f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_l}} \alpha_j^{(l)}(\mathbf{x}) \mathbb{1}_{\{k=i\}}$$

# Terms in B

Let's calculate the gradients w.r.t. last layer weights

$$\partial_{W_{ij}^{(l)}} f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_l}} \alpha_j^{(l)}(\mathbf{x}) \mathbb{1}_{\{k=i\}}$$

B

# Terms in B

Let's calculate the gradients w.r.t. last layer weights

$$\partial_{W_{ij}^{(l)}} f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_l}} \alpha_j^{(l)}(\mathbf{x}) \mathbb{1}_{\{k=i\}}$$

$$B = \sum_{i,j=1}^{d_l, d_{l+1}} \partial_{W_{ij}^{(l)}} f_{\theta,k}(\mathbf{x}) \partial_{W_{ij}^{(l)}} f_{\theta,k'}(\mathbf{x}')$$

## Terms in B

Let's calculate the gradients w.r.t. last layer weights

$$\partial_{W_{ij}^{(l)}} f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_l}} \alpha_j^{(l)}(\mathbf{x}) \mathbb{1}_{\{k=i\}}$$

$$B = \sum_{i,j=1}^{d_l, d_{l+1}} \partial_{W_{ij}^{(l)}} f_{\theta,k}(\mathbf{x}) \partial_{W_{ij}^{(l)}} f_{\theta,k'}(\mathbf{x}') = \frac{1}{d_l} \sum_{i,j=1}^{d_l, d_{l+1}} \alpha_j^{(l)}(\mathbf{x}) \mathbb{1}_{\{k=i\}} \alpha_j(\mathbf{x}') \mathbb{1}_{\{k'=i\}}$$

# Terms in B

Let's calculate the gradients w.r.t. last layer weights

$$\partial_{W_{ij}^{(l)}} f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_l}} \alpha_j^{(l)}(\mathbf{x}) \mathbb{1}_{\{k=i\}}$$

$$\begin{aligned} B &= \sum_{i,j=1}^{d_l, d_{l+1}} \partial_{W_{ij}^{(l)}} f_{\theta,k}(\mathbf{x}) \partial_{W_{ij}^{(l)}} f_{\theta,k'}(\mathbf{x}') = \frac{1}{d_l} \sum_{i,j=1}^{d_l, d_{l+1}} \alpha_j^{(l)}(\mathbf{x}) \mathbb{1}_{\{k=i\}} \alpha_j(\mathbf{x}') \mathbb{1}_{\{k'=i\}} \\ &= \frac{1}{d_l} \sum_{j=1}^{d_{l+1}} \alpha_j^{(l)}(\mathbf{x}) \alpha_j^{(l)}(\mathbf{x}') \mathbb{1}_{\{k=k'\}} \end{aligned}$$

# Terms in B

Let's calculate the gradients w.r.t. last layer weights

$$\partial_{W_{ij}^{(l)}} f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_l}} \alpha_j^{(l)}(\mathbf{x}) \mathbb{1}_{\{k=i\}}$$

$$\begin{aligned} B &= \sum_{i,j=1}^{d_l, d_{l+1}} \partial_{W_{ij}^{(l)}} f_{\theta,k}(\mathbf{x}) \partial_{W_{ij}^{(l)}} f_{\theta,k'}(\mathbf{x}') = \frac{1}{d_l} \sum_{i,j=1}^{d_l, d_{l+1}} \alpha_j^{(l)}(\mathbf{x}) \mathbb{1}_{\{k=i\}} \alpha_j(\mathbf{x}') \mathbb{1}_{\{k'=i\}} \\ &= \frac{1}{d_l} \sum_{j=1}^{d_{l+1}} \alpha_j^{(l)}(\mathbf{x}) \alpha_j^{(l)}(\mathbf{x}') \mathbb{1}_{\{k=k'\}} \xrightarrow{d_l \rightarrow \infty} \mathbb{E} \left[ \sigma(\tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x})) \sigma(\tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}')) \right] \mathbb{1}_{\{k=k'\}} \end{aligned}$$

# Terms in B

Let's calculate the gradients w.r.t. last layer weights

$$\partial_{W_{ij}^{(l)}} f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_l}} \alpha_j^{(l)}(\mathbf{x}) \mathbb{1}_{\{k=i\}}$$

$$\begin{aligned} B &= \sum_{i,j=1}^{d_l, d_{l+1}} \partial_{W_{ij}^{(l)}} f_{\theta,k}(\mathbf{x}) \partial_{W_{ij}^{(l)}} f_{\theta,k'}(\mathbf{x}') = \frac{1}{d_l} \sum_{i,j=1}^{d_l, d_{l+1}} \alpha_j^{(l)}(\mathbf{x}) \mathbb{1}_{\{k=i\}} \alpha_j(\mathbf{x}') \mathbb{1}_{\{k'=i\}} \\ &= \frac{1}{d_l} \sum_{j=1}^{d_{l+1}} \alpha_j^{(l)}(\mathbf{x}) \alpha_j^{(l)}(\mathbf{x}') \mathbb{1}_{\{k=k'\}} \xrightarrow{d_l \rightarrow \infty} \mathbb{E} \left[ \sigma(\tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x})) \sigma(\tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}')) \right] \mathbb{1}_{\{k=k'\}} \\ &= \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \tilde{\Sigma}^{(l)})} \left[ \sigma(z_1) \sigma(z_2) \right] \mathbb{1}_{\{k=k'\}} \end{aligned}$$

# Terms in B

Let's calculate the gradients w.r.t. last layer weights

$$\partial_{W_{ij}^{(l)}} f_{\theta,k}(\mathbf{x}) = \frac{1}{\sqrt{d_l}} \alpha_j^{(l)}(\mathbf{x}) \mathbb{1}_{\{k=i\}}$$

$$\begin{aligned} B &= \sum_{i,j=1}^{d_l, d_{l+1}} \partial_{W_{ij}^{(l)}} f_{\theta,k}(\mathbf{x}) \partial_{W_{ij}^{(l)}} f_{\theta,k'}(\mathbf{x}') = \frac{1}{d_l} \sum_{i,j=1}^{d_l, d_{l+1}} \alpha_j^{(l)}(\mathbf{x}) \mathbb{1}_{\{k=i\}} \alpha_j(\mathbf{x}') \mathbb{1}_{\{k'=i\}} \\ &= \frac{1}{d_l} \sum_{j=1}^{d_{l+1}} \alpha_j^{(l)}(\mathbf{x}) \alpha_j^{(l)}(\mathbf{x}') \mathbb{1}_{\{k=k'\}} \xrightarrow{d_l \rightarrow \infty} \mathbb{E} \left[ \sigma(\tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x})) \sigma(\tilde{\alpha}_{j,\infty}^{(l)}(\mathbf{x}')) \right] \mathbb{1}_{\{k=k'\}} \\ &= \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \tilde{\Sigma}^{(l)})} \left[ \sigma(z_1) \sigma(z_2) \right] \mathbb{1}_{\{k=k'\}} = \Sigma^{(l+1)}(\mathbf{x}, \mathbf{x}') \mathbb{1}_{\{k=k'\}} \end{aligned}$$

**All in All**

## All in All

- Summing up the two quantities gives:

# All in All

- Summing up the two quantities gives:

$$\Theta_{kk'}^{(l+1)}(\mathbf{x}, \mathbf{x}') = \Theta_\infty^{(l)}(\mathbf{x}, \mathbf{x}') \mathbb{1}_{\{k=k'\}} \dot{\Sigma}^{(l+1)}(\mathbf{x}, \mathbf{x}') + \Sigma^{(l+1)}(\mathbf{x}, \mathbf{x}') \mathbb{1}_{\{k=k'\}}$$

## All in All

- Summing up the two quantities gives:

$$\Theta_{kk'}^{(l+1)}(\mathbf{x}, \mathbf{x}') = \Theta_\infty^{(l)}(\mathbf{x}, \mathbf{x}') \mathbb{1}_{\{k=k'\}} \dot{\Sigma}^{(l+1)}(\mathbf{x}, \mathbf{x}') + \Sigma^{(l+1)}(\mathbf{x}, \mathbf{x}') \mathbb{1}_{\{k=k'\}}$$

- This concludes the proof!

# Time-independence

**Theorem 4:**

# Time-independence

Let us state the time-independence theorem for completeness:

**Theorem 4:**

# Time-independence

Let us state the time-independence theorem for completeness:

## Theorem 4:

Assume  $\sigma$  is Lipschitz, twice differentiable and has bounded second derivative. Fix some time horizon  $T$  and assume that

$$\int_0^T \sqrt{\sum_{i=1}^n \frac{\partial L(f(\mathbf{x}), y_i)}{\partial f(\mathbf{x})}} \Big|_{\mathbf{x}=\mathbf{x}_i} dt \text{ is bounded.}$$

# Time-independence

Let us state the time-independence theorem for completeness:

## Theorem 4:

Assume  $\sigma$  is Lipschitz, twice differentiable and has bounded second derivative. Fix some time horizon  $T$  and assume that

$$\int_0^T \sqrt{\sum_{i=1}^n \frac{\partial L(f(\mathbf{x}), y_i)}{\partial f(\mathbf{x})}} \Big|_{\mathbf{x}=\mathbf{x}_i} dt \text{ is bounded.}$$

Then:  $\Theta_t^{(L)} \xrightarrow{d_1, \dots, d_{L-1} \rightarrow \infty} \Theta_\infty^{(L)} \otimes \mathbb{1}_{d_L \times d_L}$  for all  $t \in [0, T]$

# **Proof Sketch**

# Proof Sketch

- Again relies on **induction** over the depth:

# Proof Sketch

- Again relies on **induction** over the depth:

$$\Theta_{\infty}^{(l+1)}(x, x')$$

# Proof Sketch

- Again relies on **induction** over the depth:

$$\begin{aligned}\Theta_{\infty}^{(I+1)}(\mathbf{x}, \mathbf{x}') &= \lim_{d_I \rightarrow \infty} \frac{1}{d_I} \sum_{j=1}^{d_I} W_{kj}^{(I)} W_{k'j}^{(I)} \dot{\sigma} \left( \tilde{\alpha}_j^{(I)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_j^{(I)}(\mathbf{x}') \right) \Theta_{\infty}^{(I)}(\mathbf{x}, \mathbf{x}') \\ &\quad + \frac{1}{d_I} \sum_{j=1}^{d_{I+1}} \alpha_j^{(I)}(\mathbf{x}) \alpha_j^{(I)}(\mathbf{x}') \mathbb{1}_{\{k=k'\}}\end{aligned}$$

# Proof Sketch

- Again relies on **induction** over the depth:

$$\begin{aligned}\Theta_{\infty}^{(l+1)}(\mathbf{x}, \mathbf{x}') &= \lim_{d_l \rightarrow \infty} \frac{1}{d_l} \sum_{j=1}^{d_l} W_{kj}^{(l)} W_{k'j}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_j^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_j^{(l)}(\mathbf{x}') \right) \Theta_{\infty}^{(l)}(\mathbf{x}, \mathbf{x}') \\ &\quad + \frac{1}{d_l} \sum_{j=1}^{d_{l+1}} \alpha_j^{(l)}(\mathbf{x}) \alpha_j^{(l)}(\mathbf{x}') \mathbb{1}_{\{k=k'\}}\end{aligned}$$

- The time derivative  $\partial_t \Theta_{\infty}^{(l+1)}(\mathbf{x}, \mathbf{x}')$  asymptotically looks like:

# Proof Sketch

- Again relies on **induction** over the depth:

$$\begin{aligned}\Theta_{\infty}^{(l+1)}(\mathbf{x}, \mathbf{x}') &= \lim_{d_l \rightarrow \infty} \frac{1}{d_l} \sum_{j=1}^{d_l} W_{kj}^{(l)} W_{k'j}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_j^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_j^{(l)}(\mathbf{x}') \right) \Theta_{\infty}^{(l)}(\mathbf{x}, \mathbf{x}') \\ &\quad + \frac{1}{d_l} \sum_{j=1}^{d_{l+1}} \alpha_j^{(l)}(\mathbf{x}) \alpha_j^{(l)}(\mathbf{x}') \mathbb{1}_{\{k=k'\}}\end{aligned}$$

- The time derivative  $\partial_t \Theta_{\infty}^{(l+1)}(\mathbf{x}, \mathbf{x}')$  asymptotically looks like:

$$\partial_t \Theta_{\infty}^{(l+1)}(\mathbf{x}, \mathbf{x}')$$

# Proof Sketch

- Again relies on **induction** over the depth:

$$\begin{aligned}\Theta_{\infty}^{(l+1)}(\mathbf{x}, \mathbf{x}') &= \lim_{d_l \rightarrow \infty} \frac{1}{d_l} \sum_{j=1}^{d_l} W_{kj}^{(l)} W_{k'j}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_j^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_j^{(l)}(\mathbf{x}') \right) \Theta_{\infty}^{(l)}(\mathbf{x}, \mathbf{x}') \\ &\quad + \frac{1}{d_l} \sum_{j=1}^{d_{l+1}} \alpha_j^{(l)}(\mathbf{x}) \alpha_j^{(l)}(\mathbf{x}') \mathbb{1}_{\{k=k'\}}\end{aligned}$$

- The time derivative  $\partial_t \Theta_{\infty}^{(l+1)}(\mathbf{x}, \mathbf{x}')$  asymptotically looks like:

$$\partial_t \Theta_{\infty}^{(l+1)}(\mathbf{x}, \mathbf{x}') \sim \lim_{d_l \rightarrow \infty} \frac{1}{d_l} \sum_{j=1}^{d_l} \mathcal{O}(\partial_t W_{kj}^{(l)}) + \mathcal{O}(\partial_t \tilde{\alpha}_j^{(l)}) + \mathcal{O}(\partial_t \alpha_j^{(l)})$$

# Proof Sketch

- Again relies on **induction** over the depth:

$$\begin{aligned}\Theta_{\infty}^{(l+1)}(\mathbf{x}, \mathbf{x}') &= \lim_{d_l \rightarrow \infty} \frac{1}{d_l} \sum_{j=1}^{d_l} W_{kj}^{(l)} W_{k'j}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_j^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_j^{(l)}(\mathbf{x}') \right) \Theta_{\infty}^{(l)}(\mathbf{x}, \mathbf{x}') \\ &\quad + \frac{1}{d_l} \sum_{j=1}^{d_{l+1}} \alpha_j^{(l)}(\mathbf{x}) \alpha_j^{(l)}(\mathbf{x}') \mathbb{1}_{\{k=k'\}}\end{aligned}$$

- The time derivative  $\partial_t \Theta_{\infty}^{(l+1)}(\mathbf{x}, \mathbf{x}')$  asymptotically looks like:

$$\begin{aligned}\partial_t \Theta_{\infty}^{(l+1)}(\mathbf{x}, \mathbf{x}') &\sim \lim_{d_l \rightarrow \infty} \frac{1}{d_l} \sum_{j=1}^{d_l} \mathcal{O}(\partial_t W_{kj}^{(l)}) + \mathcal{O}(\partial_t \tilde{\alpha}_j^{(l)}) + \mathcal{O}(\partial_t \alpha_j^{(l)}) \\ &\sim \mathcal{O}(\partial_t W_{kj}^{(l)}) + \mathcal{O}(\partial_t \tilde{\alpha}_j^{(l)}) + \mathcal{O}(\partial_t \alpha_j^{(l)})\end{aligned}$$

# Proof Sketch

- Again relies on **induction** over the depth:

$$\begin{aligned}\Theta_{\infty}^{(l+1)}(\mathbf{x}, \mathbf{x}') &= \lim_{d_l \rightarrow \infty} \frac{1}{d_l} \sum_{j=1}^{d_l} W_{kj}^{(l)} W_{k'j}^{(l)} \dot{\sigma} \left( \tilde{\alpha}_j^{(l)}(\mathbf{x}) \right) \dot{\sigma} \left( \tilde{\alpha}_j^{(l)}(\mathbf{x}') \right) \Theta_{\infty}^{(l)}(\mathbf{x}, \mathbf{x}') \\ &\quad + \frac{1}{d_l} \sum_{j=1}^{d_{l+1}} \alpha_j^{(l)}(\mathbf{x}) \alpha_j^{(l)}(\mathbf{x}') \mathbb{1}_{\{k=k'\}}\end{aligned}$$

- The time derivative  $\partial_t \Theta_{\infty}^{(l+1)}(\mathbf{x}, \mathbf{x}')$  asymptotically looks like:

$$\begin{aligned}\partial_t \Theta_{\infty}^{(l+1)}(\mathbf{x}, \mathbf{x}') &\sim \lim_{d_l \rightarrow \infty} \frac{1}{d_l} \sum_{j=1}^{d_l} \mathcal{O}(\partial_t W_{kj}^{(l)}) + \mathcal{O}(\partial_t \tilde{\alpha}_j^{(l)}) + \mathcal{O}(\partial_t \alpha_j^{(l)}) \\ &\sim \mathcal{O}(\partial_t W_{kj}^{(l)}) + \mathcal{O}(\partial_t \tilde{\alpha}_j^{(l)}) + \mathcal{O}(\partial_t \alpha_j^{(l)}) = \mathcal{O}\left(\frac{1}{\sqrt{d_l}}\right)\end{aligned}$$

# **NTK Animation**

# NTK Animation

- Train on just two images of *MNIST*  $\implies \Theta^{(L)}(\mathbf{X}, \mathbf{X}) \in \mathbb{R}^{2 \times 2}$

# NTK Animation

- Train on just two images of *MNIST*  $\implies \Theta^{(L)}(\mathbf{X}, \mathbf{X}) \in \mathbb{R}^{2 \times 2}$
- Plot corresponding ellipse through training for different widths

# NTK Animation

- Train on just two images of *MNIST*  $\implies \Theta^{(L)}(\mathbf{X}, \mathbf{X}) \in \mathbb{R}^{2 \times 2}$
- Plot corresponding ellipse through training for different widths

# **Discussion**

# Discussion

- Training NNs  $\approx$  kernel regression with NTK  $\implies$  The corresponding kernel literature can be connected to NNs concerning generalization, optimization etc

# Discussion

- Training NNs  $\approx$  kernel regression with NTK  $\implies$  The corresponding kernel literature can be connected to NNs concerning generalization, optimization etc
- Only gradient flow, Seyed will present a paper that also talks about **gradient descent**

# Discussion

- Training NNs  $\approx$  kernel regression with NTK  $\implies$  The corresponding kernel literature can be connected to NNs concerning generalization, optimization etc
- Only gradient flow, Seyed will present a paper that also talks about **gradient descent**
- Concerning a generalization bound on loss:

# Discussion

- Training NNs  $\approx$  kernel regression with NTK  $\implies$  The corresponding kernel literature can be connected to NNs concerning generalization, optimization etc
- Only gradient flow, Seyed will present a paper that also talks about **gradient descent**
- Concerning a generalization bound on loss:

Bound  $\geq$  observed kernel performance  $\geq$  observed NN performance

# Discussion

- Training NNs  $\approx$  kernel regression with NTK  $\implies$  The corresponding kernel literature can be connected to NNs concerning generalization, optimization etc
- Only gradient flow, Seyed will present a paper that also talks about **gradient descent**
- Concerning a generalization bound on loss:

Bound  $\geq$  observed kernel performance  $\geq$  observed NN performance

- Can hence probably not explain why neural networks work **so** well if we take the kernel route but it may be a good start

# Discussion

- Training NNs  $\approx$  kernel regression with NTK  $\implies$  The corresponding kernel literature can be connected to NNs concerning generalization, optimization etc
- Only gradient flow, Seyed will present a paper that also talks about **gradient descent**
- Concerning a generalization bound on loss:

Bound  $\geq$  observed kernel performance  $\geq$  observed NN performance

- Can hence probably not explain why neural networks work **so** well if we take the kernel route but it may be a good start
- We observe that NNs work better with increasing width. This is a bit conflicting with infinite width as there performance decreases (Aurelien will present a paper that shows this empirically)

# **Backup Slides**

# Closed Form Solution (I)

## Closed Form Solution (I)

- Plugging this gradient into our ODE gives:

## Closed Form Solution (I)

- Plugging this gradient into our ODE gives:

$$\partial_t f_t^\infty(\mathbf{x}) = - \left( \Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T (f_t^\infty(\mathbf{X}) - \mathbf{y})$$

## Closed Form Solution (I)

- Plugging this gradient into our ODE gives:

$$\partial_t f_t^\infty(\mathbf{x}) = - \left( \Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T (f_t^\infty(\mathbf{X}) - \mathbf{y})$$

That's a system of **linear first order** differential equations!

## Closed Form Solution (I)

- Plugging this gradient into our ODE gives:

$$\partial_t f_t^\infty(\mathbf{x}) = - \left( \Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T (f_t^\infty(\mathbf{X}) - \mathbf{y})$$

That's a system of **linear first order** differential equations!

- Careful, every  $\mathbf{x}$  induces a **different** function  $f_t^\infty(\mathbf{x})$  in  $t$ !

# Closed Form Solution (I)

- Plugging this gradient into our ODE gives:

$$\partial_t f_t^\infty(\mathbf{x}) = - \left( \Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X}) \right)^T (f_t^\infty(\mathbf{X}) - \mathbf{y})$$

That's a system of **linear first order** differential equations!

- Careful, every  $\mathbf{x}$  induces a **different** function  $f_t^\infty(\mathbf{x})$  in  $t$ !
- We first need to solve it on the training set  $\mathbf{X}$ :

$$\partial_t f_t^\infty(\mathbf{X}) = -\Theta_\infty(\mathbf{X}, \mathbf{X})(f_t^\infty(\mathbf{X}) - \mathbf{y})$$

## Closed Form Solution (II)

## Closed Form Solution (II)

- Recall:  $\dot{x}(t) = A(x(t) - x^*) \implies x(t) = x^* + e^{At}(x(0) - x^*)$

## Closed Form Solution (II)

- Recall:  $\dot{\mathbf{x}}(t) = A(\mathbf{x}(t) - \mathbf{x}^*) \implies \mathbf{x}(t) = \mathbf{x}^* + e^{At}(\mathbf{x}(0) - \mathbf{x}^*)$
- We hence get the following solution on the data

$$f_t^\infty(\mathbf{X}) = \mathbf{y} + e^{-t\Theta_\infty^{(L)}(\mathbf{X}, \mathbf{X})}(f_0^\infty(\mathbf{X}) - \mathbf{y})$$

## Closed Form Solution (II)

- Recall:  $\dot{\mathbf{x}}(t) = A(\mathbf{x}(t) - \mathbf{x}^*) \implies \mathbf{x}(t) = \mathbf{x}^* + e^{At}(\mathbf{x}(0) - \mathbf{x}^*)$
- We hence get the following solution on the data

$$f_t^\infty(\mathbf{X}) = \mathbf{y} + e^{-t\Theta_\infty^{(L)}(\mathbf{X}, \mathbf{X})}(f_0^\infty(\mathbf{X}) - \mathbf{y})$$

- Getting the function value at an arbitrary point  $\mathbf{x}$  is now just a matter of **integrating** the original ODE

## Closed Form Solution (II)

- Recall:  $\dot{\mathbf{x}}(t) = A(\mathbf{x}(t) - \mathbf{x}^*) \implies \mathbf{x}(t) = \mathbf{x}^* + e^{At}(\mathbf{x}(0) - \mathbf{x}^*)$
- We hence get the following solution on the data

$$f_t^\infty(\mathbf{X}) = \mathbf{y} + e^{-t\Theta_\infty^{(L)}(\mathbf{X}, \mathbf{X})}(f_0^\infty(\mathbf{X}) - \mathbf{y})$$

- Getting the function value at an arbitrary point  $\mathbf{x}$  is now just a matter of **integrating** the original ODE
- Some more tricks:

## Closed Form Solution (II)

- Recall:  $\dot{\mathbf{x}}(t) = A(\mathbf{x}(t) - \mathbf{x}^*) \implies \mathbf{x}(t) = \mathbf{x}^* + e^{At}(\mathbf{x}(0) - \mathbf{x}^*)$
- We hence get the following solution on the data

$$f_t^\infty(\mathbf{X}) = \mathbf{y} + e^{-t\Theta_\infty^{(L)}(\mathbf{X}, \mathbf{X})}(f_0^\infty(\mathbf{X}) - \mathbf{y})$$

- Getting the function value at an arbitrary point  $\mathbf{x}$  is now just a matter of **integrating** the original ODE
- Some more tricks:
  - 1)  $\int \mathbf{a}^T \mathbf{X}(t) \mathbf{b} dt = \mathbf{a}^T (\int \mathbf{X}(t) dt) \mathbf{b}$

## Closed Form Solution (II)

- Recall:  $\dot{\mathbf{x}}(t) = A(\mathbf{x}(t) - \mathbf{x}^*) \implies \mathbf{x}(t) = \mathbf{x}^* + e^{At}(\mathbf{x}(0) - \mathbf{x}^*)$
- We hence get the following solution on the data

$$f_t^\infty(\mathbf{X}) = \mathbf{y} + e^{-t\Theta_\infty^{(L)}(\mathbf{X}, \mathbf{X})} (f_0^\infty(\mathbf{X}) - \mathbf{y})$$

- Getting the function value at an arbitrary point  $\mathbf{x}$  is now just a matter of **integrating** the original ODE
- Some more tricks:

$$1) \int \mathbf{a}^T \mathbf{X}(t) \mathbf{b} dt = \mathbf{a}^T (\int \mathbf{X}(t) dt) \mathbf{b}$$

$$2) \int e^{At} dt = A^{-1} (e^{At} - \mathbb{1}) + C$$

## Closed Form Solution (III)

Performing this integral gives

$$\begin{aligned} f_t^\infty(\mathbf{x}) &= \int -\left(\Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X})\right)^T (f_t^\infty(\mathbf{X}) - \mathbf{y}) \ dt \\ &= \int -\left(\Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X})\right)^T e^{-\Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X})t} (f_0^\infty(\mathbf{X}) - \mathbf{y}) \ dt \\ &\stackrel{1)}{=} -\left(\Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X})\right)^T \left( \int e^{-\Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X})t} \ dt \right) (f_0^\infty(\mathbf{X}) - \mathbf{y}) \\ &\stackrel{2)}{=} \left(\Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X})\right)^T \left( \left(\Theta_\infty^{(L)}\right)^{-1} (e^{-\Theta_\infty^{(L)}(\mathbf{x}, \mathbf{X})t} - 1) \right) (f_0^\infty(\mathbf{X}) - \mathbf{y}) \\ &\quad + f_0^\infty(\mathbf{x}) \end{aligned}$$