

Logistics

Logistics

- Reading Group on **Neural Tangent Kernel**

Logistics

- Reading Group on **Neural Tangent Kernel**
- For me, **generalization** view point but very interested in other insights from optimization etc

Logistics

- Reading Group on **Neural Tangent Kernel**
- For me, **generalization** view point but very interested in other insights from optimization etc
- Weekly meetings every **Wednesday afternoon**, time flexible but "default" at **2-3pm**

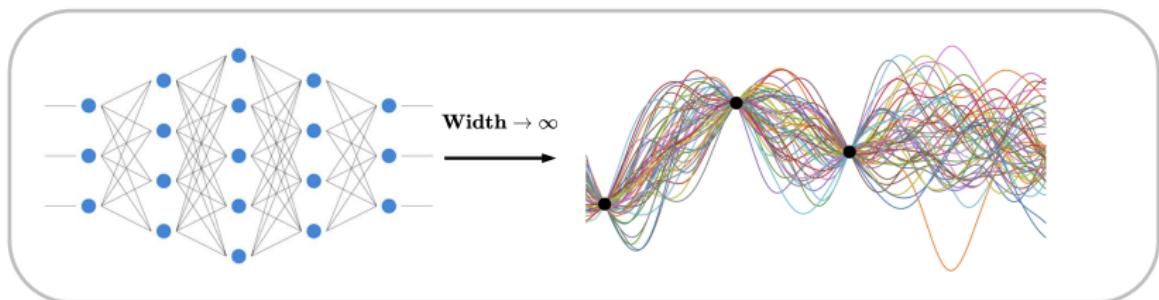
Logistics

- Reading Group on **Neural Tangent Kernel**
- For me, **generalization** view point but very interested in other insights from optimization etc
- Weekly meetings every **Wednesday afternoon**, time flexible but "default" at **2-3pm**
- Sign up for a paper or add your own paper of interest

Logistics

- Reading Group on **Neural Tangent Kernel**
- For me, **generalization** view point but very interested in other insights from optimization etc
- Weekly meetings every **Wednesday afternoon**, time flexible but "default" at **2-3pm**
- Sign up for a paper or add your own paper of interest
- **Virtual** meetings via Zoom (unless of course everyone is back at the lab) and hence slide presentations

Deep Neural Networks as Gaussian Processes



**Jaehoon Lee, Yasaman Bahri, Roman Novak,
Samuel S. Schoenholz, Jeffrey Pennington, Jascha
Sohl-Dickstein**

Recap: Gaussian Process

Recap: Gaussian Process

- Mathematical definition of a Gaussian process:

Recap: Gaussian Process

- Mathematical definition of a Gaussian process:

Definition: A **Gaussian process** $\{X_t\}_{t=0}^{\infty} \sim \mathcal{GP}(0, \Sigma)$ is a stochastic process such that for any finite collection $(t_1, \dots, t_n) \in \mathbb{R}_+^n$, we have $(X_{t_1}, \dots, X_{t_n}) \sim \mathcal{N}(\mathbf{0}, \Sigma)$ where $\Sigma_{ij} = \Sigma(X_{t_i}, X_{t_j})$ and $\Sigma : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$

Recap: Gaussian Process

- Mathematical definition of a Gaussian process:

Definition: A **Gaussian process** $\{X_t\}_{t=0}^{\infty} \sim \mathcal{GP}(0, \Sigma)$ is a stochastic process such that for any finite collection $(t_1, \dots, t_n) \in \mathbb{R}_+^n$, we have $(X_{t_1}, \dots, X_{t_n}) \sim \mathcal{N}(\mathbf{0}, \Sigma)$ where $\Sigma_{ij} = \Sigma(X_{t_i}, X_{t_j})$ and $\Sigma : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$

- Allows to sample random functions where **path continuity** depends solely on Σ (*Kolmogorov's continuity theorem*)

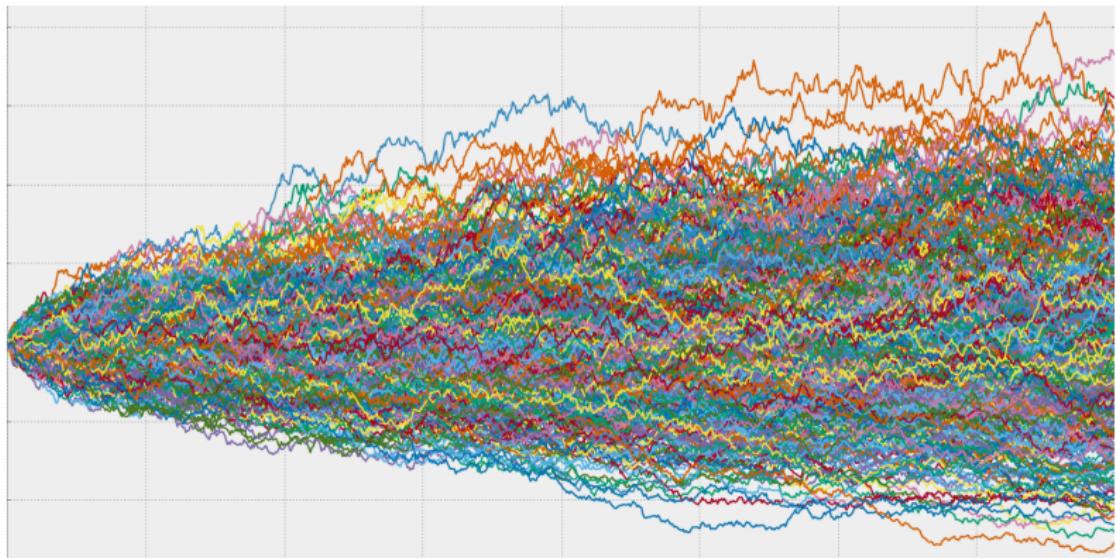
Recap: Gaussian Process

- Mathematical definition of a Gaussian process:

Definition: A **Gaussian process** $\{X_t\}_{t=0}^{\infty} \sim \mathcal{GP}(0, \Sigma)$ is a stochastic process such that for any finite collection $(t_1, \dots, t_n) \in \mathbb{R}_+^n$, we have $(X_{t_1}, \dots, X_{t_n}) \sim \mathcal{N}(\mathbf{0}, \Sigma)$ where $\Sigma_{ij} = \Sigma(X_{t_i}, X_{t_j})$ and $\Sigma : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$

- Allows to sample random functions where **path continuity** depends solely on Σ (*Kolmogorov's continuity theorem*)
- **Example:** Brownian motion with $\Sigma(s, t) = \min(s, t)$

Sample Paths of Brownian Motion



Gaussian Processes in ML (I)

Gaussian Processes in ML (I)

- **Regression Problem:** Observations $\mathbf{x}_1, \dots, \mathbf{x}_n \subset \mathbb{R}^d$ and a noisy model $y_i = f(\mathbf{x}_i) + \epsilon_i$ where $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

Gaussian Processes in ML (I)

- **Regression Problem:** Observations $\mathbf{x}_1, \dots, \mathbf{x}_n \subset \mathbb{R}^d$ and a noisy model $y_i = f(\mathbf{x}_i) + \epsilon_i$ where $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$
- **Prior over functions:** $f(\cdot) \sim \mathcal{GP}(0, \Sigma)$

But how to model Σ ?

Gaussian Processes in ML (I)

- **Regression Problem:** Observations $\mathbf{x}_1, \dots, \mathbf{x}_n \subset \mathbb{R}^d$ and a noisy model $y_i = f(\mathbf{x}_i) + \epsilon_i$ where $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

- **Prior over functions:** $f(\cdot) \sim \mathcal{GP}(0, \Sigma)$

But how to model Σ ?

- Use information in the features \mathbf{x}_i by taking approach

$$\Sigma_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$$

where $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is a **symmetric** and **positive semi-definite** kernel

Gaussian Processes in ML (II)

Gaussian Processes in ML (II)

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

Gaussian Processes in ML (II)

- $[K(\mathbf{x}, \mathbf{X})]_i = K(\mathbf{x}, \mathbf{x}_i)$ and $[K(\mathbf{X}, \mathbf{X})]_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$
- **Bayes:** Given new observation $\mathbf{x}^* \in \mathbb{R}^d$, how to **update** prior to get $P(y^* | \mathbf{y}, \mathbf{X}, \mathbf{x}^*)$?

Gaussian Processes in ML (II)

- $[K(\mathbf{x}, \mathbf{X})]_i = K(\mathbf{x}, \mathbf{x}_i)$ and $[K(\mathbf{X}, \mathbf{X})]_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$
- **Bayes:** Given new observation $\mathbf{x}^* \in \mathbb{R}^d$, how to **update** prior to get $P(y^*|\mathbf{y}, \mathbf{X}, \mathbf{x}^*)$?
- First obtain joint distribution:

Gaussian Processes in ML (II)

- $[K(\mathbf{x}, \mathbf{X})]_i = K(\mathbf{x}, \mathbf{x}_i)$ and $[K(\mathbf{X}, \mathbf{X})]_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$
- **Bayes:** Given new observation $\mathbf{x}^* \in \mathbb{R}^d$, how to **update** prior to get $P(y^* | \mathbf{y}, \mathbf{X}, \mathbf{x}^*)$?
- First obtain joint distribution:

$$\begin{pmatrix} \mathbf{y} \\ y^* \end{pmatrix} = \underbrace{\begin{pmatrix} h(\mathbf{X}) \\ h(\mathbf{x}^*) \end{pmatrix}}_{+} + \underbrace{\begin{pmatrix} \boldsymbol{\epsilon} \\ \epsilon^* \end{pmatrix}}$$

Gaussian Processes in ML (II)

- $[K(\mathbf{x}, \mathbf{X})]_i = K(\mathbf{x}, \mathbf{x}_i)$ and $[K(\mathbf{X}, \mathbf{X})]_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$
- **Bayes:** Given new observation $\mathbf{x}^* \in \mathbb{R}^d$, how to **update** prior to get $P(y^* | \mathbf{y}, \mathbf{X}, \mathbf{x}^*)$?
- First obtain joint distribution:

$$\begin{pmatrix} \mathbf{y} \\ y^* \end{pmatrix} = \underbrace{\begin{pmatrix} h(\mathbf{X}) \\ h(\mathbf{x}^*) \end{pmatrix}}_{\sim \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} K(\mathbf{X}, \mathbf{X}) & K(\mathbf{X}, \mathbf{x}) \\ K(\mathbf{X}, \mathbf{x}) & K(\mathbf{x}, \mathbf{x}) \end{pmatrix}\right)} + \underbrace{\begin{pmatrix} \boldsymbol{\epsilon} \\ \epsilon^* \end{pmatrix}}$$

Gaussian Processes in ML (II)

- $[K(\mathbf{x}, \mathbf{X})]_i = K(\mathbf{x}, \mathbf{x}_i)$ and $[K(\mathbf{X}, \mathbf{X})]_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$
- **Bayes:** Given new observation $\mathbf{x}^* \in \mathbb{R}^d$, how to **update** prior to get $P(y^* | \mathbf{y}, \mathbf{X}, \mathbf{x}^*)$?
- First obtain joint distribution:

$$\begin{pmatrix} \mathbf{y} \\ y^* \end{pmatrix} = \underbrace{\begin{pmatrix} h(\mathbf{X}) \\ h(\mathbf{x}^*) \end{pmatrix}}_{\sim \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} K(\mathbf{X}, \mathbf{X}) & K(\mathbf{x}, \mathbf{X}) \\ K(\mathbf{X}, \mathbf{x}) & K(\mathbf{x}, \mathbf{x}) \end{pmatrix}\right)} + \underbrace{\begin{pmatrix} \boldsymbol{\epsilon} \\ \epsilon^* \end{pmatrix}}_{\sim \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} \sigma^2 \mathbb{1} & 0 \\ 0 & \sigma^2 \end{pmatrix}\right)}$$

Gaussian Processes in ML (III)

Gaussian Processes in ML (III)

Adding the two distributions gives

Gaussian Processes in ML (III)

Adding the two distributions gives

$$\begin{pmatrix} \mathbf{y} \\ y^* \end{pmatrix} \Big| (\mathbf{X}, \mathbf{x}^*) \sim \mathcal{N} \left(\mathbf{0}, \begin{pmatrix} K(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbb{1} & K(\mathbf{X}, \mathbf{x}) \\ K(\mathbf{X}, \mathbf{x}) & K(\mathbf{x}, \mathbf{x}) + \sigma^2 \end{pmatrix} \right)$$

Gaussian Processes in ML (III)

Adding the two distributions gives

$$\begin{pmatrix} \mathbf{y} \\ y^* \end{pmatrix} \Big| (\mathbf{X}, \mathbf{x}^*) \sim \mathcal{N} \left(\mathbf{0}, \begin{pmatrix} K(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbb{1} & K(\mathbf{x}, \mathbf{X}) \\ K(\mathbf{X}, \mathbf{x}) & K(\mathbf{x}, \mathbf{x}) + \sigma^2 \end{pmatrix} \right)$$

- **Fact:** For $\begin{pmatrix} z_A \\ z_B \end{pmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{pmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{pmatrix} \right)$ we have

$$x_A | x_B \sim \mathcal{N} \left(\Sigma_{AB} \Sigma_{BB}^{-1} x_B, \Sigma_{AA} - \Sigma_{AB} \Sigma_{BB}^{-1} \Sigma_{BA} \right)$$

Gaussian Processes in ML (III)

Adding the two distributions gives

$$\begin{pmatrix} \mathbf{y} \\ y^* \end{pmatrix} \Big| (\mathbf{X}, \mathbf{x}^*) \sim \mathcal{N} \left(\mathbf{0}, \begin{pmatrix} K(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbb{1} & K(\mathbf{x}, \mathbf{X}) \\ K(\mathbf{X}, \mathbf{x}) & K(\mathbf{x}, \mathbf{x}) + \sigma^2 \end{pmatrix} \right)$$

- **Fact:** For $\begin{pmatrix} \mathbf{z}_A \\ \mathbf{z}_B \end{pmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{pmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{pmatrix} \right)$ we have

$$\mathbf{x}_A | \mathbf{x}_B \sim \mathcal{N} \left(\Sigma_{AB} \Sigma_{BB}^{-1} \mathbf{x}_B, \Sigma_{AA} - \Sigma_{AB} \Sigma_{BB}^{-1} \Sigma_{BA} \right)$$

- We hence have $y^* | (\mathbf{X}, \mathbf{y}, \mathbf{x}^*) \sim \mathcal{N}(\mu^*, \sigma^*)$ where

$$1 \quad \mu^* = K(\mathbf{x}^*, \mathbf{X})^T (K(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbb{1})^{-1} \mathbf{y}$$

$$2 \quad \sigma^* = K(\mathbf{x}^*, \mathbf{x}^*) - K(\mathbf{x}^*, \mathbf{X})^T (K(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbb{1})^{-1} K(\mathbf{x}^*, \mathbf{X})$$

Gaussian Processes in ML (IV)

Gaussian Processes in ML (IV)

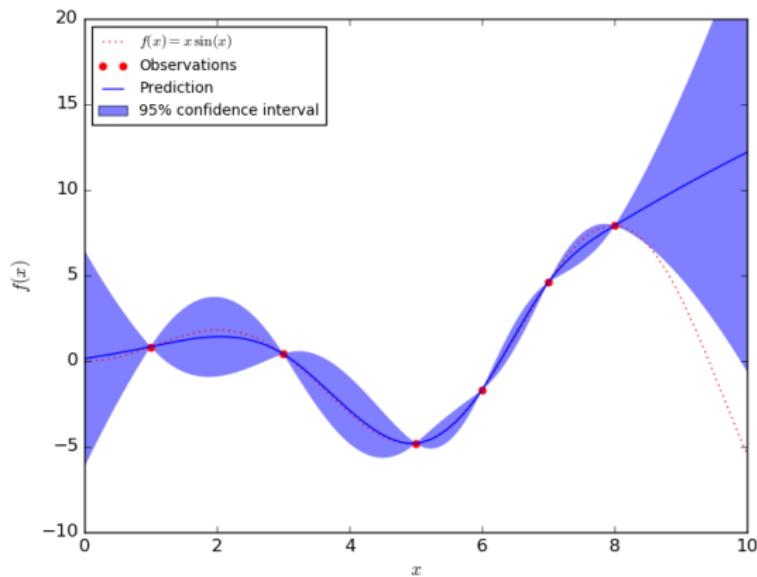
- We have a **probability distribution** as a prediction for y^* and hence a measure for uncertainty

Gaussian Processes in ML (IV)

- We have a **probability distribution** as a prediction for y^* and hence a measure for uncertainty
- Can take the **mean** μ^* or sample from f

Gaussian Processes in ML (IV)

- We have a **probability distribution** as a prediction for y^* and hence a measure for uncertainty
- Can take the **mean** μ^* or sample from f



Neural Networks: Notation

Neural Networks: Notation

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

Neural Networks: Notation

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

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

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

Neural Networks: 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 $\mathbf{W}^{(l)} \in \mathbb{R}^{d_{l+1} \times d_l}$, $\mathbf{b}^{(l)} \in \mathbb{R}^{d_{l+1}}$, σ 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

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

Multidimensional Central Limit Theorem

Multidimensional Central Limit Theorem

- CLT will be major **ingredient** for the proofs

Multidimensional Central Limit Theorem

- CLT will be major **ingredient** for the proofs
- Let us restate it for completeness:

Multidimensional Central Limit Theorem

- CLT will be major **ingredient** for the proofs
- Let us restate it for completeness:

Theorem: For $\mathbf{X}_1 = \begin{pmatrix} X_{11} \\ \vdots \\ X_{1d} \end{pmatrix}, \dots, \mathbf{X}_n = \begin{pmatrix} X_{n1} \\ \vdots \\ X_{nd} \end{pmatrix}$ i.i.d. multi-dim.

random variables with $\mathbb{E}[\mathbf{X}_i] = \mathbf{0}$ and $\text{cov}(\mathbf{X}_i) = \Sigma$, it holds that in the limit

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{X}_i \xrightarrow{(d)} \mathcal{N}(\mathbf{0}, \Sigma)$$

Main Result

Main Result

Let's state the **main result** of this paper:

Main Result

Let's state the **main result** of this paper:

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

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

Main Result

Let's state the **main result** of this paper:

Theorem: *Given a NN f_θ , 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

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

Main Result

Let's state the **main result** of this paper:

Theorem: *Given a NN f_θ , 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

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

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

Some Comments

Some Comments

- Consider the **infinite-width** limit $d_1, \dots, d_{L-1} \rightarrow \infty$ (all limits here denote convergence in law)

Some Comments

- Consider the **infinite-width** limit $d_1, \dots, d_{L-1} \rightarrow \infty$ (all limits here denote convergence in **law**)
- **In words:** Each component of the output converges to a Gaussian process with a known covariance function modulo an integral over σ (known for *ReLU*)

Some Comments

- Consider the **infinite-width** limit $d_1, \dots, d_{L-1} \rightarrow \infty$ (all limits here denote convergence in **law**)
- **In words:** Each component of the output converges to a Gaussian process with a known covariance function modulo an integral over σ (known for *ReLU*)
- Allows to perform full **Bayesian-inference** with an infinitely wide neural network (no Monte Carlo methods)

Some Comments

- Consider the **infinite-width** limit $d_1, \dots, d_{L-1} \rightarrow \infty$ (all limits here denote convergence in **law**)
- **In words:** Each component of the output converges to a Gaussian process with a known covariance function modulo an integral over σ (known for *ReLU*)
- Allows to perform full **Bayesian-inference** with an infinitely wide neural network (no Monte Carlo methods)
- Only study network at **initialization**, no training at all

Some Comments

- Consider the **infinite-width** limit $d_1, \dots, d_{L-1} \rightarrow \infty$ (all limits here denote convergence in **law**)
- **In words:** Each component of the output converges to a Gaussian process with a known covariance function modulo an integral over σ (known for *ReLU*)
- Allows to perform full **Bayesian-inference** with an infinitely wide neural network (no Monte Carlo methods)
- Only study network at **initialization**, no training at all
- Theory would actually break down since weights become **correlated** after just one step of SGD

Proof of Theorem

Proof of Theorem

- Proof uses **induction** on the depth L

Proof of Theorem

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

Proof of Theorem

- Proof uses **induction** on the depth L
- Limits $d_1, \dots, d_{L-1} \rightarrow \infty$ are taken **sequentially**
- Central Limit Theorem performs all the heavy-lifting

Proof of Theorem

- Proof uses **induction** on the depth L
- Limits $d_1, \dots, d_{L-1} \rightarrow \infty$ are taken **sequentially**
- Central Limit Theorem performs all the heavy-lifting
- Use $\Sigma^{(l)}$ for covariance **function** and $\Sigma^{(l)}$ for the resulting covariance **matrix** for data x_1, \dots, x_k

Proof of Theorem

- Proof uses **induction** on the depth L
- Limits $d_1, \dots, d_{L-1} \rightarrow \infty$ are taken **sequentially**
- Central Limit Theorem performs all the heavy-lifting
- Use $\Sigma^{(l)}$ for covariance **function** and $\Sigma^{(l)}$ for the resulting covariance **matrix** for data x_1, \dots, x_k
- Write limits as f_θ^∞ and $\tilde{\alpha}_{i,\infty}^{(l)}$

Base Case: $L = 1$

Base Case: $L = 1$

Network takes the form

$$f_{\theta}(\mathbf{x}) = \frac{1}{\sqrt{d_0}} \mathbf{W}^{(0)} \mathbf{x} + \beta \mathbf{b}^{(0)}$$

Base Case: $L = 1$

Network takes the form

$$f_{\theta}(\mathbf{x}) = \frac{1}{\sqrt{d_0}} \mathbf{W}^{(0)} \mathbf{x} + \beta \mathbf{b}^{(0)}$$


← Remains fixed

Base Case: $L = 1$

Network takes the form

$$f_{\theta}(\mathbf{x}) = \frac{1}{\sqrt{d_0}} \mathbf{W}^{(0)} \mathbf{x} + \beta \mathbf{b}^{(0)}$$


← Remains fixed

Considering just one component leads to

Base Case: $L = 1$

Network takes the form

$$f_{\theta}(\mathbf{x}) = \frac{1}{\sqrt{d_0}} \mathbf{W}^{(0)} \mathbf{x} + \beta \mathbf{b}^{(0)}$$

← Remains fixed

Considering just one component leads to

$$f_{\theta,i}(\mathbf{x}) = \frac{1}{\sqrt{d_0}} \overbrace{\sum_{j=1}^{d_0} W_{ij}^{(0)} x_j}^{} + \beta b_i^{(0)}$$

Base Case: $L = 1$

Network takes the form

$$f_{\theta}(\mathbf{x}) = \frac{1}{\sqrt{d_0}} \mathbf{W}^{(0)} \mathbf{x} + \beta \mathbf{b}^{(0)}$$

 Remains fixed

Considering just one component leads to

$$f_{\theta,i}(\mathbf{x}) = \frac{1}{\sqrt{d_0}} \underbrace{\sum_{j=1}^{d_0} W_{ij}^{(0)} x_j}_{\text{Lin.comb.of i.i.d. } \mathcal{N}} + \beta b_i^{(0)}$$



Base Case: $L = 1$

Network takes the form

$$f_{\theta}(\mathbf{x}) = \frac{1}{\sqrt{d_0}} \mathbf{W}^{(0)} \mathbf{x} + \beta \mathbf{b}^{(0)}$$

 Remains fixed

Considering just one component leads to

$$\begin{aligned} f_{\theta,i}(\mathbf{x}) &= \frac{1}{\sqrt{d_0}} \underbrace{\sum_{j=1}^{d_0} W_{ij}^{(0)} x_j}_{\text{Lin.comb.of i.i.d. } \mathcal{N}} + \beta b_i^{(0)} \\ &\sim \mathcal{N} \left(0, \frac{1}{d_0} \sum_{j=1}^{d_0} \text{Var} \left(W_{ij}^{(0)} \right) x_j^2 + \beta^2 \text{Var}(b_i^{(0)}) \right) \end{aligned}$$

Base Case: $L = 1$

Network takes the form

$$f_{\theta}(\mathbf{x}) = \frac{1}{\sqrt{d_0}} \mathbf{W}^{(0)} \mathbf{x} + \beta \mathbf{b}^{(0)}$$

 Remains fixed

Considering just one component leads to

$$\begin{aligned} f_{\theta,i}(\mathbf{x}) &= \frac{1}{\sqrt{d_0}} \underbrace{\sum_{j=1}^{d_0} W_{ij}^{(0)} x_j}_{\text{Lin.comb.of i.i.d. } \mathcal{N}} + \beta b_i^{(0)} \\ &\sim \mathcal{N} \left(0, \frac{1}{d_0} \sum_{j=1}^{d_0} \text{Var} \left(W_{ij}^{(0)} \right) x_j^2 + \beta^2 \text{Var}(b_i^{(0)}) \right) \\ &\stackrel{(d)}{=} \mathcal{N} \left(0, \frac{1}{d_0} \|\mathbf{x}\|_2^2 + \beta^2 \right) \end{aligned}$$

Covariance Structure

Covariance Structure

For any $x, x' \in \mathbb{R}^{d_0}$ we have:

Covariance Structure

For any $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{d_0}$ we have:

$$\Sigma^{(1)}(\mathbf{x}, \mathbf{x}') = \mathbb{E}[f_{\theta,i}(\mathbf{x}), f_{\theta,i}(\mathbf{x}')] \quad \text{for } i = 1, \dots, n$$

Covariance Structure

For any $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{d_0}$ we have:

$$\Sigma^{(1)}(\mathbf{x}, \mathbf{x}') = \mathbb{E}[f_{\theta,i}(\mathbf{x}), f_{\theta,i}(\mathbf{x}')]$$

$$= \mathbb{E}\left[\left(\frac{1}{\sqrt{d_0}} \sum_{j=1}^{d_0} W_{ij}^{(0)} x_j + \beta b_i^{(0)}\right) \left(\frac{1}{\sqrt{d_0}} \sum_{j=1}^{d_0} W_{ij}^{(0)} x'_j + \beta b_i^{(0)}\right)\right]$$

Covariance Structure

For any $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{d_0}$ we have:

$$\Sigma^{(1)}(\mathbf{x}, \mathbf{x}') = \mathbb{E}[f_{\theta,i}(\mathbf{x}), f_{\theta,i}(\mathbf{x}')] =$$

$$= \mathbb{E}\left[\left(\frac{1}{\sqrt{d_0}} \sum_{j=1}^{d_0} W_{ij}^{(0)} x_j + \beta b_i^{(0)}\right) \left(\frac{1}{\sqrt{d_0}} \sum_{j=1}^{d_0} W_{ij}^{(0)} x'_j + \beta b_i^{(0)}\right)\right]$$

$$\mathbb{E}[W_{ij}^{(0)} b_i^{(0)}] = 0$$

Covariance Structure

For any $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{d_0}$ we have:

$$\Sigma^{(1)}(\mathbf{x}, \mathbf{x}') = \mathbb{E}[f_{\theta,i}(\mathbf{x}), f_{\theta,i}(\mathbf{x}')] =$$

$$\mathbb{E}[W_{ij}^{(0)} b_i^{(0)}] = 0$$

$$= \mathbb{E}\left[\left(\frac{1}{\sqrt{d_0}} \sum_{j=1}^{d_0} W_{ij}^{(0)} x_j + \beta b_i^{(0)}\right) \left(\frac{1}{\sqrt{d_0}} \sum_{j=1}^{d_0} W_{ij}^{(0)} x'_j + \beta b_i^{(0)}\right)\right]$$

$$= \frac{1}{d_0} \sum_{j=1}^{d_0} \sum_{l=1}^{d_0} \mathbb{E}[W_{ij}^{(0)} W_{il}^{(0)}] x_j x'_l + \beta^2$$

Covariance Structure

For any $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{d_0}$ we have:

$$\Sigma^{(1)}(\mathbf{x}, \mathbf{x}') = \mathbb{E}[f_{\theta,i}(\mathbf{x}), f_{\theta,i}(\mathbf{x}')] =$$

$$\mathbb{E}[W_{ij}^{(0)} b_i^{(0)}] = 0$$

$$\begin{aligned} &= \mathbb{E}\left[\left(\frac{1}{\sqrt{d_0}} \sum_{j=1}^{d_0} W_{ij}^{(0)} x_j + \beta b_i^{(0)}\right) \left(\frac{1}{\sqrt{d_0}} \sum_{j=1}^{d_0} W_{ij}^{(0)} x'_j + \beta b_i^{(0)}\right)\right] \\ &= \frac{1}{d_0} \sum_{j=1}^{d_0} \sum_{l=1}^{d_0} \mathbb{E}[W_{ij}^{(0)} W_{il}^{(0)}] x_j x'_l + \beta^2 \\ &= \frac{1}{d_0} \sum_{j=1}^{d_0} \sum_{l=1}^{d_0} \mathbb{1}_{\{j=l\}} x_j x'_l + \beta^2 \end{aligned}$$

Covariance Structure

For any $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{d_0}$ we have:

$$\Sigma^{(1)}(\mathbf{x}, \mathbf{x}') = \mathbb{E}[f_{\theta,i}(\mathbf{x}), f_{\theta,i}(\mathbf{x}')] =$$

$$\mathbb{E}[W_{ij}^{(0)} b_i^{(0)}] = 0$$

$$= \mathbb{E}\left[\left(\frac{1}{\sqrt{d_0}} \sum_{j=1}^{d_0} W_{ij}^{(0)} x_j + \beta b_i^{(0)}\right) \left(\frac{1}{\sqrt{d_0}} \sum_{j=1}^{d_0} W_{ij}^{(0)} x'_j + \beta b_i^{(0)}\right)\right]$$

$$= \frac{1}{d_0} \sum_{j=1}^{d_0} \sum_{l=1}^{d_0} \mathbb{E}[W_{ij}^{(0)} W_{il}^{(0)}] x_j x'_l + \beta^2$$

$$= \frac{1}{d_0} \sum_{j=1}^{d_0} \sum_{l=1}^{d_0} \mathbb{1}_{\{j=l\}} x_j x'_l + \beta^2$$

$$\mathbb{E}[W_{ij}^{(0)} W_{il}^{(0)}] = \mathbb{1}_{\{j=l\}}$$

Covariance Structure

For any $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{d_0}$ we have:

$$\Sigma^{(1)}(\mathbf{x}, \mathbf{x}') = \mathbb{E}[f_{\theta,i}(\mathbf{x}), f_{\theta,i}(\mathbf{x}')] =$$

$$\mathbb{E}[W_{ij}^{(0)} b_i^{(0)}] = 0$$

$$\begin{aligned} &= \mathbb{E}\left[\left(\frac{1}{\sqrt{d_0}} \sum_{j=1}^{d_0} W_{ij}^{(0)} x_j + \beta b_i^{(0)}\right) \left(\frac{1}{\sqrt{d_0}} \sum_{j=1}^{d_0} W_{ij}^{(0)} x'_j + \beta b_i^{(0)}\right)\right] \\ &= \frac{1}{d_0} \sum_{j=1}^{d_0} \sum_{l=1}^{d_0} \mathbb{E}[W_{ij}^{(0)} W_{il}^{(0)}] x_j x'_l + \beta^2 \\ &= \frac{1}{d_0} \sum_{j=1}^{d_0} \sum_{l=1}^{d_0} \mathbb{1}_{\{j=l\}} x_j x'_l + \beta^2 \\ &= \frac{1}{d_0} \mathbf{x}^T \mathbf{x}' + \beta^2 \end{aligned}$$

$$\mathbb{E}[W_{ij}^{(0)} W_{il}^{(0)}] = \mathbb{1}_{\{j=l\}}$$

Gaussian Process Property

Gaussian Process Property

- Moreover, for $i \neq j$ and all \mathbf{x}, \mathbf{x}' we have:

Gaussian Process Property

- Moreover, for $i \neq j$ and all \mathbf{x}, \mathbf{x}' we have:

$$\mathbb{E}[f_{\theta,i}(\mathbf{x}), f_{\theta,j}(\mathbf{x}')] = 0$$

Gaussian Process Property

- Moreover, for $i \neq j$ and all \mathbf{x}, \mathbf{x}' we have:

$$\mathbb{E}[W_{ik}^{(0)} W_{jk}^{(0)}] = 0$$
$$\mathbb{E}[f_{\theta,i}(\mathbf{x}), f_{\theta,j}(\mathbf{x}')] = 0$$

Gaussian Process Property

- Moreover, for $i \neq j$ and all \mathbf{x}, \mathbf{x}' we have:

$$\mathbb{E}[W_{ik}^{(0)} W_{jk}^{(0)}] = 0$$
$$\mathbb{E}[f_{\theta,i}(\mathbf{x}), f_{\theta,j}(\mathbf{x}')] = 0$$

which implies **independence** due to Gaussianity.

Gaussian Process Property

- Moreover, for $i \neq j$ and all \mathbf{x}, \mathbf{x}' we have:

$$\mathbb{E}[W_{ik}^{(0)} W_{jk}^{(0)}] = 0$$
$$\mathbb{E}[f_{\theta,i}(\mathbf{x}), f_{\theta,j}(\mathbf{x}')] = 0$$

which implies **independence** due to Gaussianity.

- Finally, we have for any $\mathbf{x}^1, \dots, \mathbf{x}^k \in \mathbb{R}^{d_0}$:

Gaussian Process Property

- Moreover, for $i \neq j$ and all \mathbf{x}, \mathbf{x}' we have:

$$\mathbb{E}[W_{ik}^{(0)} W_{jk}^{(0)}] = 0$$
$$\mathbb{E}[f_{\theta,i}(\mathbf{x}), f_{\theta,j}(\mathbf{x}')] = 0$$

which implies **independence** due to Gaussianity.

- Finally, we have for any $\mathbf{x}^1, \dots, \mathbf{x}^k \in \mathbb{R}^{d_0}$:

$$\left(f_{\theta,i}(\mathbf{x}^1), \dots, f_{\theta,i}(\mathbf{x}^k) \right) \sim \mathcal{N} \left(\mathbf{0}, \Sigma^{(1)} \right)$$

Gaussian Process Property

- Moreover, for $i \neq j$ and all \mathbf{x}, \mathbf{x}' we have:

$$\mathbb{E}[W_{ik}^{(0)} W_{jk}^{(0)}] = 0$$
$$\mathbb{E}[f_{\theta,i}(\mathbf{x}), f_{\theta,j}(\mathbf{x}')] = 0$$

which implies **independence** due to Gaussianity.

- Finally, we have for any $\mathbf{x}^1, \dots, \mathbf{x}^k \in \mathbb{R}^{d_0}$:

$$\left(f_{\theta,i}(\mathbf{x}^1), \dots, f_{\theta,i}(\mathbf{x}^k) \right) \sim \mathcal{N} \left(\mathbf{0}, \Sigma^{(1)} \right)$$

where again $\Sigma_{ij}^{(1)} = \Sigma^{(1)}(\mathbf{x}_i, \mathbf{x}_j)$

- All in all:** $f_{\theta,i} \sim \mathcal{GP}(0, \Sigma^{(1)})$

Induction Step

Induction Step

- **Assume:** $\tilde{\alpha}_{i,\infty}^{(l)}(\cdot) \sim \mathcal{GP}(0, \Sigma^{(l)})$ for $d_1, \dots, d_{l-1} \rightarrow \infty$

Induction Step

- **Assume:** $\tilde{\alpha}_{i,\infty}^{(l)}(\cdot) \sim \mathcal{GP}(0, \Sigma^{(l)})$ for $d_1, \dots, d_{l-1} \rightarrow \infty$
- Consider network of depth $l + 1$:

Induction Step

- **Assume:** $\tilde{\alpha}_{i,\infty}^{(l)}(\cdot) \sim \mathcal{GP}(0, \Sigma^{(l)})$ for $d_1, \dots, d_{l-1} \rightarrow \infty$
- Consider network of depth $l + 1$:

$$f_{\theta}(\mathbf{x}) = \tilde{\alpha}^{(l+1)}(\mathbf{x}) = \frac{1}{\sqrt{d_l}} \mathbf{W}^{(l)} \alpha^{(l)}(\mathbf{x}) + \beta \mathbf{b}^{(l)}$$

Induction Step

- **Assume:** $\tilde{\alpha}_{i,\infty}^{(l)}(\cdot) \sim \mathcal{GP}(0, \Sigma^{(l)})$ for $d_1, \dots, d_{l-1} \rightarrow \infty$
- Consider network of depth $l + 1$:

$$f_{\theta}(\mathbf{x}) = \tilde{\alpha}^{(l+1)}(\mathbf{x}) = \frac{1}{\sqrt{d_l}} \mathbf{W}^{(l)} \alpha^{(l)}(\mathbf{x}) + \beta \mathbf{b}^{(l)}$$

- Component-wise:

Induction Step

- **Assume:** $\tilde{\alpha}_{i,\infty}^{(l)}(\cdot) \sim \mathcal{GP}(0, \Sigma^{(l)})$ for $d_1, \dots, d_{l-1} \rightarrow \infty$
- Consider network of depth $l+1$:

$$f_{\theta}(\mathbf{x}) = \tilde{\alpha}^{(l+1)}(\mathbf{x}) = \frac{1}{\sqrt{d_l}} \mathbf{W}^{(l)} \alpha^{(l)}(\mathbf{x}) + \beta \mathbf{b}^{(l)}$$

- Component-wise:

$$f_{\theta,i}(\mathbf{x}) = \frac{1}{\sqrt{d_l}} \sum_{k=1}^{d_l} W_{ik}^{(l)} \alpha_k^{(l)}(\mathbf{x}) + \beta b_i^{(l)}$$

Induction Step

Induction Step

- Take $(\mathbf{x}^1, \dots, \mathbf{x}^k) \subset \mathbb{R}^{d_0}$ and look at $(f_{\theta,i}(\mathbf{x}^1), \dots, f_{\theta,i}(\mathbf{x}^k))$

Induction Step

- Take $(\mathbf{x}^1, \dots, \mathbf{x}^k) \subset \mathbb{R}^{d_0}$ and look at $(f_{\theta,i}(\mathbf{x}^1), \dots, f_{\theta,i}(\mathbf{x}^k))$
- We know that for $d_1, \dots, d_{l-1} \rightarrow \infty$, $f_{\theta,i}(\mathbf{x})$ is a sum of i.i.d. random variables:

Induction Step

- Take $(\mathbf{x}^1, \dots, \mathbf{x}^k) \subset \mathbb{R}^{d_0}$ and look at $(f_{\theta,i}(\mathbf{x}^1), \dots, f_{\theta,i}(\mathbf{x}^k))$
- We know that for $d_1, \dots, d_{l-1} \rightarrow \infty$, $f_{\theta,i}(\mathbf{x})$ is a sum of i.i.d. random variables:

$$f_{\theta,i}(\mathbf{x}) \xrightarrow{d_1, \dots, d_{l-1} \rightarrow \infty} \frac{1}{\sqrt{d_l}} \sum_{k=1}^{d_l} W_{ik}^{(l)} \alpha_{k,\infty}^{(l)}(\mathbf{x}) + \beta b_i^{(l)}$$

Induction Step

- Take $(\mathbf{x}^1, \dots, \mathbf{x}^k) \subset \mathbb{R}^{d_0}$ and look at $(f_{\theta,i}(\mathbf{x}^1), \dots, f_{\theta,i}(\mathbf{x}^k))$
- We know that for $d_1, \dots, d_{l-1} \rightarrow \infty$, $f_{\theta,i}(\mathbf{x})$ is a sum of i.i.d. random variables:

$$f_{\theta,i}(\mathbf{x}) \xrightarrow{d_1, \dots, d_{l-1} \rightarrow \infty} \frac{1}{\sqrt{d_l}} \sum_{k=1}^{d_l} W_{ik}^{(l)} \alpha_{k,\infty}^{(l)}(\mathbf{x}) + \beta b_i^{(l)}$$

- The covariance structure needed for CLT can be written as

Induction Step

- Take $(\mathbf{x}^1, \dots, \mathbf{x}^k) \subset \mathbb{R}^{d_0}$ and look at $(f_{\theta,i}(\mathbf{x}^1), \dots, f_{\theta,i}(\mathbf{x}^k))$
- We know that for $d_1, \dots, d_{I-1} \rightarrow \infty$, $f_{\theta,i}(\mathbf{x})$ is a sum of i.i.d. random variables:

$$f_{\theta,i}(\mathbf{x}) \xrightarrow{d_1, \dots, d_{I-1} \rightarrow \infty} \frac{1}{\sqrt{d_I}} \sum_{k=1}^{d_I} W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}) + \beta b_i^{(I)}$$

- The covariance structure needed for CLT can be written as

$$\hat{\Sigma}_{st} = \text{cov} \left(W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^s), W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^t) \right)$$

Induction Step

- Take $(\mathbf{x}^1, \dots, \mathbf{x}^k) \subset \mathbb{R}^{d_0}$ and look at $(f_{\theta,i}(\mathbf{x}^1), \dots, f_{\theta,i}(\mathbf{x}^k))$
- We know that for $d_1, \dots, d_{I-1} \rightarrow \infty$, $f_{\theta,i}(\mathbf{x})$ is a sum of i.i.d. random variables:

$$f_{\theta,i}(\mathbf{x}) \xrightarrow{d_1, \dots, d_{I-1} \rightarrow \infty} \frac{1}{\sqrt{d_I}} \sum_{k=1}^{d_I} W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}) + \beta b_i^{(I)}$$

- The covariance structure needed for CLT can be written as

$$\begin{aligned}\hat{\Sigma}_{st} &= \text{cov} \left(W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^s), W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^t) \right) \\ &= \mathbb{E}[W_{ik}^{(L)} W_{ik}^{(I)}] \mathbb{E}[\alpha_{k,\infty}^{(I)}(\mathbf{x}^s) \alpha_{k,\infty}^{(I)}(\mathbf{x}^t)]\end{aligned}$$

Induction Step

- Take $(\mathbf{x}^1, \dots, \mathbf{x}^k) \subset \mathbb{R}^{d_0}$ and look at $(f_{\theta,i}(\mathbf{x}^1), \dots, f_{\theta,i}(\mathbf{x}^k))$
- We know that for $d_1, \dots, d_{l-1} \rightarrow \infty$, $f_{\theta,i}(\mathbf{x})$ is a sum of i.i.d. random variables:

$$f_{\theta,i}(\mathbf{x}) \xrightarrow{d_1, \dots, d_{l-1} \rightarrow \infty} \frac{1}{\sqrt{d_l}} \sum_{k=1}^{d_l} W_{ik}^{(l)} \alpha_{k,\infty}^{(l)}(\mathbf{x}) + \beta b_i^{(l)}$$

- The covariance structure needed for CLT can be written as

$$\begin{aligned}\hat{\Sigma}_{st} &= \text{cov} \left(W_{ik}^{(l)} \alpha_{k,\infty}^{(l)}(\mathbf{x}^s), W_{ik}^{(l)} \alpha_{k,\infty}^{(l)}(\mathbf{x}^t) \right) \\ &= \mathbb{E}[W_{ik}^{(l)} W_{ik}^{(l)}] \mathbb{E}[\alpha_{k,\infty}^{(l)}(\mathbf{x}^s) \alpha_{k,\infty}^{(l)}(\mathbf{x}^t)] \\ &= \mathbb{E}[\alpha_{k,\infty}^{(l)}(\mathbf{x}^s) \alpha_{k,\infty}^{(l)}(\mathbf{x}^t)]\end{aligned}$$

Covariance Structure

Covariance Structure

- By the multidimensional CLT we now obtain

Covariance Structure

- By the multidimensional CLT we now obtain

$$\begin{pmatrix} f_{\theta,i}(\mathbf{x}^1) \\ \vdots \\ f_{\theta,i}(\mathbf{x}^k) \end{pmatrix}$$

Covariance Structure

- By the multidimensional CLT we now obtain

$$\begin{pmatrix} f_{\theta,i}(\mathbf{x}^1) \\ \vdots \\ f_{\theta,i}(\mathbf{x}^k) \end{pmatrix} = \frac{1}{\sqrt{d_I}} \sum_{k=1}^{d_I} \begin{pmatrix} W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^1) \\ \vdots \\ W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^k) \end{pmatrix}$$

Covariance Structure

- By the multidimensional CLT we now obtain

$$\begin{pmatrix} f_{\theta,i}(\mathbf{x}^1) \\ \vdots \\ f_{\theta,i}(\mathbf{x}^k) \end{pmatrix} = \frac{1}{\sqrt{d_I}} \sum_{k=1}^{d_I} \begin{pmatrix} W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^1) \\ \vdots \\ W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^k) \end{pmatrix} \xrightarrow{d_I \rightarrow \infty} \mathcal{N}\left(\mathbf{0}, \hat{\Sigma} + \beta^2\right)$$

Covariance Structure

- By the multidimensional CLT we now obtain

$$\begin{pmatrix} f_{\theta,i}(\mathbf{x}^1) \\ \vdots \\ f_{\theta,i}(\mathbf{x}^k) \end{pmatrix} = \frac{1}{\sqrt{d_I}} \sum_{k=1}^{d_I} \begin{pmatrix} W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^1) \\ \vdots \\ W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^k) \end{pmatrix} \xrightarrow{d_I \rightarrow \infty} \mathcal{N} \left(\mathbf{0}, \hat{\Sigma} + \beta^2 \right)$$

- Hence: $f_{\theta,i}^\infty(\cdot) \sim \mathcal{GP}(0, \Sigma^{(I+1)})$

Covariance Structure

- By the multidimensional CLT we now obtain

$$\begin{pmatrix} f_{\theta,i}(\mathbf{x}^1) \\ \vdots \\ f_{\theta,i}(\mathbf{x}^k) \end{pmatrix} = \frac{1}{\sqrt{d_I}} \sum_{k=1}^{d_I} \begin{pmatrix} W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^1) \\ \vdots \\ W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^k) \end{pmatrix} \xrightarrow{d_I \rightarrow \infty} \mathcal{N} \left(\mathbf{0}, \hat{\Sigma} + \beta^2 \right)$$

- Hence: $f_{\theta,i}^\infty(\cdot) \sim \mathcal{GP}(0, \Sigma^{(I+1)})$
- Moreover, we can write out the covariance function as

Covariance Structure

- By the multidimensional CLT we now obtain

$$\begin{pmatrix} f_{\theta,i}(\mathbf{x}^1) \\ \vdots \\ f_{\theta,i}(\mathbf{x}^k) \end{pmatrix} = \frac{1}{\sqrt{d_I}} \sum_{k=1}^{d_I} \begin{pmatrix} W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^1) \\ \vdots \\ W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^k) \end{pmatrix} \xrightarrow{d_I \rightarrow \infty} \mathcal{N}\left(\mathbf{0}, \hat{\Sigma} + \beta^2\right)$$

- Hence: $f_{\theta,i}^\infty(\cdot) \sim \mathcal{GP}(0, \Sigma^{(I+1)})$
- Moreover, we can write out the covariance function as

$$\Sigma^{(I+1)}(\mathbf{x}, \mathbf{x}') = \hat{\Sigma}(\mathbf{x}, \mathbf{x}') + \beta^2$$

Covariance Structure

- By the multidimensional CLT we now obtain

$$\begin{pmatrix} f_{\theta,i}(\mathbf{x}^1) \\ \vdots \\ f_{\theta,i}(\mathbf{x}^k) \end{pmatrix} = \frac{1}{\sqrt{d_I}} \sum_{k=1}^{d_I} \begin{pmatrix} W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^1) \\ \vdots \\ W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^k) \end{pmatrix} \xrightarrow{d_I \rightarrow \infty} \mathcal{N}\left(\mathbf{0}, \hat{\Sigma} + \beta^2\right)$$

- Hence: $f_{\theta,i}^\infty(\cdot) \sim \mathcal{GP}(0, \Sigma^{(I+1)})$
- Moreover, we can write out the covariance function as

$$\Sigma^{(I+1)}(\mathbf{x}, \mathbf{x}') = \hat{\Sigma}(\mathbf{x}, \mathbf{x}') + \beta^2 = \mathbb{E}[\alpha_{k,\infty}^{(I+1)}(\mathbf{x}) \alpha_{k,\infty}^{(I+1)}(\mathbf{x}')] + \beta^2$$

Covariance Structure

- By the multidimensional CLT we now obtain

$$\begin{pmatrix} f_{\theta,i}(\mathbf{x}^1) \\ \vdots \\ f_{\theta,i}(\mathbf{x}^k) \end{pmatrix} = \frac{1}{\sqrt{d_I}} \sum_{k=1}^{d_I} \begin{pmatrix} W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^1) \\ \vdots \\ W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^k) \end{pmatrix} \xrightarrow{d_I \rightarrow \infty} \mathcal{N} \left(\mathbf{0}, \hat{\Sigma} + \beta^2 \right)$$

- Hence: $f_{\theta,i}^\infty(\cdot) \sim \mathcal{GP}(0, \Sigma^{(I+1)})$
- Moreover, we can write out the covariance function as

$$\begin{aligned} \Sigma^{(I+1)}(\mathbf{x}, \mathbf{x}') &= \hat{\Sigma}(\mathbf{x}, \mathbf{x}') + \beta^2 = \mathbb{E}[\alpha_{k,\infty}^{(I+1)}(\mathbf{x}) \alpha_{k,\infty}^{(I+1)}(\mathbf{x}')] + \beta^2 \\ &= \mathbb{E} \left[\sigma \left(\tilde{\alpha}_{k,\infty}^{(I+1)}(\mathbf{x}) \right) \sigma \left(\tilde{\alpha}_{k,\infty}^{(I+1)}(\mathbf{x}') \right) \right] + \beta^2 \end{aligned}$$

Covariance Structure

- By the multidimensional CLT we now obtain

$$\begin{pmatrix} f_{\theta,i}(\mathbf{x}^1) \\ \vdots \\ f_{\theta,i}(\mathbf{x}^k) \end{pmatrix} = \frac{1}{\sqrt{d_I}} \sum_{k=1}^{d_I} \begin{pmatrix} W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^1) \\ \vdots \\ W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^k) \end{pmatrix} \xrightarrow{d_I \rightarrow \infty} \mathcal{N}\left(\mathbf{0}, \hat{\Sigma} + \beta^2\right)$$

- Hence: $f_{\theta,i}^\infty(\cdot) \sim \mathcal{GP}(0, \Sigma^{(I+1)})$
- Moreover, we can write out the covariance function as

$$\begin{aligned} \Sigma^{(I+1)}(\mathbf{x}, \mathbf{x}') &= \hat{\Sigma}(\mathbf{x}, \mathbf{x}') + \beta^2 = \mathbb{E}[\alpha_{k,\infty}^{(I+1)}(\mathbf{x}) \alpha_{k,\infty}^{(I+1)}(\mathbf{x}')] + \beta^2 \\ &= \mathbb{E}\left[\sigma\left(\tilde{\alpha}_{k,\infty}^{(I+1)}(\mathbf{x})\right) \sigma\left(\tilde{\alpha}_{k,\infty}^{(I+1)}(\mathbf{x}')\right)\right] + \beta^2 \\ &= \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \tilde{\Sigma}^{(I)})} [\sigma(z_1) \sigma(z_2)] + \beta^2 \end{aligned}$$

Covariance Structure

- By the multidimensional CLT we now obtain

$$\begin{pmatrix} f_{\theta,i}(\mathbf{x}^1) \\ \vdots \\ f_{\theta,i}(\mathbf{x}^k) \end{pmatrix} = \frac{1}{\sqrt{d_I}} \sum_{k=1}^{d_I} \begin{pmatrix} W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^1) \\ \vdots \\ W_{ik}^{(I)} \alpha_{k,\infty}^{(I)}(\mathbf{x}^k) \end{pmatrix} \xrightarrow{d_I \rightarrow \infty} \mathcal{N}\left(\mathbf{0}, \hat{\Sigma} + \beta^2\right)$$

- Hence: $f_{\theta,i}^\infty(\cdot) \sim \mathcal{GP}(0, \Sigma^{(I+1)})$
- Moreover, we can write out the covariance function as

$$\begin{aligned} \Sigma^{(I+1)}(\mathbf{x}, \mathbf{x}') &= \hat{\Sigma}(\mathbf{x}, \mathbf{x}') + \beta^2 = \mathbb{E}[\alpha_{k,\infty}^{(I+1)}(\mathbf{x}) \alpha_{k,\infty}^{(I+1)}(\mathbf{x}')] + \beta^2 \\ &= \mathbb{E}\left[\sigma\left(\tilde{\alpha}_{k,\infty}^{(I+1)}(\mathbf{x})\right) \sigma\left(\tilde{\alpha}_{k,\infty}^{(I+1)}(\mathbf{x}')\right)\right] + \beta^2 \\ &= \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \tilde{\Sigma}^{(I)})} [\sigma(z_1) \sigma(z_2)] + \beta^2 \end{aligned}$$

where $\tilde{\Sigma}^{(I)} = \begin{pmatrix} \Sigma^{(I)}(\mathbf{x}, \mathbf{x}) & \Sigma^{(I)}(\mathbf{x}, \mathbf{x}') \\ \Sigma^{(I)}(\mathbf{x}', \mathbf{x}) & \Sigma^{(I)}(\mathbf{x}', \mathbf{x}') \end{pmatrix}$

Independence

Independence

- Finally we have to check **independence**:

Independence

- Finally we have to check **independence**:

$$Q_{ij} = \text{cov}(f_{\theta,i}(\mathbf{x}), f_{\theta,j}(\mathbf{x}))$$

Independence

- Finally we have to check **independence**:

$$Q_{ij} = \text{cov}(f_{\theta,i}(\mathbf{x}), f_{\theta,j}(\mathbf{x})) = \mathbb{E}[f_{\theta,i}(\mathbf{x})f_{\theta,j}(\mathbf{x})]$$

Independence

- Finally we have to check **independence**:

$$\begin{aligned} Q_{ij} &= \text{cov}(f_{\theta,i}(\mathbf{x}), f_{\theta,j}(\mathbf{x})) = \mathbb{E}[f_{\theta,i}(\mathbf{x})f_{\theta,j}(\mathbf{x})] \\ &= \frac{1}{d_I} \sum_{k=1}^{d_I} \sum_{s=1}^{d_I} \mathbb{E}[W_{ik}^{(I)} W_{js}^{(I)}] \mathbb{E}[\alpha_k^{(I)}(\mathbf{x}) \alpha_s^{(I)}(\mathbf{x})] + 0 \end{aligned}$$

Independence

- Finally we have to check **independence**:

$$\begin{aligned} Q_{ij} &= \text{cov}(f_{\theta,i}(\mathbf{x}), f_{\theta,j}(\mathbf{x})) = \mathbb{E}[f_{\theta,i}(\mathbf{x})f_{\theta,j}(\mathbf{x})] \\ &= \frac{1}{d_I} \sum_{k=1}^{d_I} \sum_{s=1}^{d_I} \mathbb{E}[W_{ik}^{(I)} W_{js}^{(I)}] \mathbb{E}[\alpha_k^{(I)}(\mathbf{x}) \alpha_s^{(I)}(\mathbf{x})] + 0 \\ &= 0 \end{aligned}$$

Independence

- Finally we have to check **independence**:

$$\begin{aligned} Q_{ij} &= \text{cov}(f_{\theta,i}(\mathbf{x}), f_{\theta,j}(\mathbf{x})) = \mathbb{E}[f_{\theta,i}(\mathbf{x})f_{\theta,j}(\mathbf{x})] \\ &= \frac{1}{d_I} \sum_{k=1}^{d_I} \sum_{s=1}^{d_I} \mathbb{E}[W_{ik}^{(I)} W_{js}^{(I)}] \mathbb{E}[\alpha_k^{(I)}(\mathbf{x}) \alpha_s^{(I)}(\mathbf{x})] + 0 \\ &= 0 \end{aligned}$$

- Again by CLT, the two variables $f_{\theta,i}(\mathbf{x})$ and $f_{\theta,j}(\mathbf{x})$ converge **jointly** to the limit $\mathcal{N}(\mathbf{0}, \mathbf{Q})$.

Independence

- Finally we have to check **independence**:

$$\begin{aligned} Q_{ij} &= \text{cov}(f_{\theta,i}(\mathbf{x}), f_{\theta,j}(\mathbf{x})) = \mathbb{E}[f_{\theta,i}(\mathbf{x})f_{\theta,j}(\mathbf{x})] \\ &= \frac{1}{d_I} \sum_{k=1}^{d_I} \sum_{s=1}^{d_I} \mathbb{E}[W_{ik}^{(I)} W_{js}^{(I)}] \mathbb{E}[\alpha_k^{(I)}(\mathbf{x}) \alpha_s^{(I)}(\mathbf{x})] + 0 \\ &= 0 \end{aligned}$$

- Again by CLT, the two variables $f_{\theta,i}(\mathbf{x})$ and $f_{\theta,j}(\mathbf{x})$ converge **jointly** to the limit $\mathcal{N}(\mathbf{0}, \mathbf{Q})$.
- Thus due to Gaussianity, $Q_{ij} = 0$ also implies independence in the limit.

Independence

- Finally we have to check **independence**:

$$\begin{aligned} Q_{ij} &= \text{cov}(f_{\theta,i}(\mathbf{x}), f_{\theta,j}(\mathbf{x})) = \mathbb{E}[f_{\theta,i}(\mathbf{x})f_{\theta,j}(\mathbf{x})] \\ &= \frac{1}{d_I} \sum_{k=1}^{d_I} \sum_{s=1}^{d_I} \mathbb{E}[W_{ik}^{(I)} W_{js}^{(I)}] \mathbb{E}[\alpha_k^{(I)}(\mathbf{x}) \alpha_s^{(I)}(\mathbf{x})] + 0 \\ &= 0 \end{aligned}$$

- Again by CLT, the two variables $f_{\theta,i}(\mathbf{x})$ and $f_{\theta,j}(\mathbf{x})$ converge **jointly** to the limit $\mathcal{N}(\mathbf{0}, \mathbf{Q})$.
- Thus due to Gaussianity, $Q_{ij} = 0$ also implies independence in the limit.
- This **concludes** the proof.

Some Empirical Illustrations

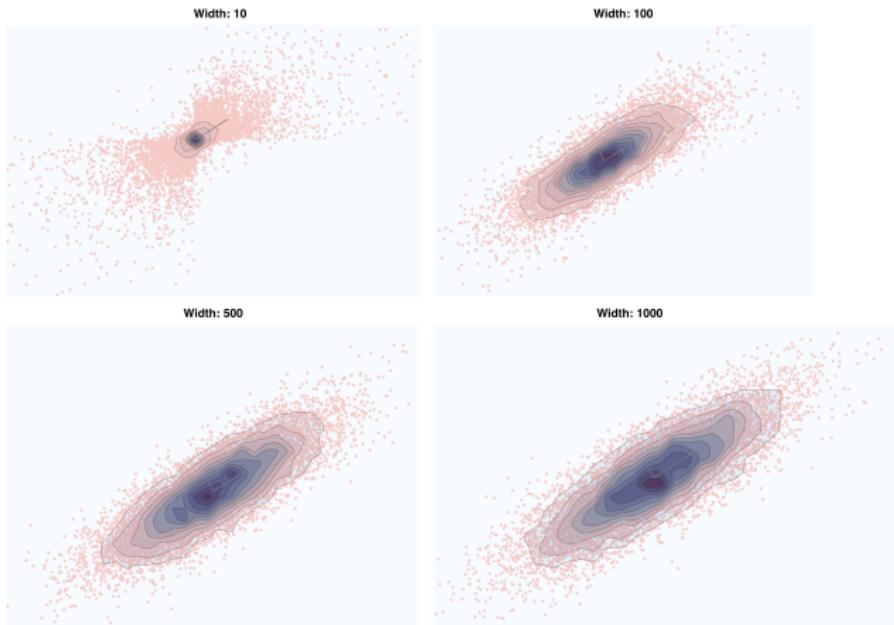
- 8-layer neural network with fixed width d on two inputs x, x' from *MNIST*

Some Empirical Illustrations

- 8-layer neural network with fixed width d on two inputs \mathbf{x}, \mathbf{x}' from *MNIST*
- Propagated the inputs through 10000 random networks and plotted $f_{\theta,1}(\mathbf{x})$ vs. $f_{\theta,1}(\mathbf{x}')$ for $d \in \{10, 100, 500, 1000\}$

Some Empirical Illustrations

- 8-layer neural network with fixed width d on two inputs \mathbf{x}, \mathbf{x}' from *MNIST*
- Propagated the inputs through 10000 random networks and plotted $f_{\theta,1}(\mathbf{x})$ vs. $f_{\theta,1}(\mathbf{x}')$ for $d \in \{10, 100, 500, 1000\}$



Most Important Takeaway (Regarding NTK)

- We can fully characterize the distribution of $f_\theta^\infty(\cdot)$ at initialization

Most Important Takeaway (Regarding NTK)

- We can fully characterize the distribution of $f_\theta^\infty(\cdot)$ at **initialization**
- We needed a scaling of $\frac{1}{\sqrt{d_l}}$ at each layer for the **asymptotics** to work

Most Important Takeaway (Regarding NTK)

- We can fully characterize the distribution of $f_\theta^\infty(\cdot)$ at **initialization**
- We needed a scaling of $\frac{1}{\sqrt{d_l}}$ at each layer for the **asymptotics** to work
- We have a **recursive** covariance structure which will reappear for NTK:

$$\Sigma^{(l)}(x, x') = \mathbb{E}_{z \sim \mathcal{N}(0, \tilde{\Sigma}^{(l-1)})} [\sigma(z_1)\sigma(z_2)] + \beta^2$$

Most Important Takeaway (Regarding NTK)

- We can fully characterize the distribution of $f_\theta^\infty(\cdot)$ at **initialization**
- We needed a scaling of $\frac{1}{\sqrt{d_l}}$ at each layer for the **asymptotics** to work
- We have a **recursive** covariance structure which will reappear for NTK:

$$\Sigma^{(l)}(x, x') = \mathbb{E}_{z \sim \mathcal{N}(0, \tilde{\Sigma}^{(l-1)})} [\sigma(z_1)\sigma(z_2)] + \beta^2$$

- Can we describe $f_\theta^\infty(\cdot)$ during **training**?

Discussion

Discussion

- All results only hold at **initialization**

Discussion

- All results only hold at **initialization**
- Seems hard to assess how the correlation structure between weights changes over course of training

Discussion

- All results only hold at **initialization**
- Seems hard to assess how the correlation structure between weights changes over course of training
- Even if we know how correlation evolves, CLT with correlated variables are hard to apply

Discussion

- All results only hold at **initialization**
- Seems hard to assess how the correlation structure between weights changes over course of training
- Even if we know how correlation evolves, CLT with correlated variables are hard to apply
- Infinite-width limit interesting because generalization gets better with **increasing width** (empirically)