# Disentangling Trainability and Generalization in Deep Neural Networks

Lechao Xiao,    Jeffrey Pennington,    Samuel Schoenholz

# Notation

# Notation

- Fully-connected L-layer network:

$$z^{(l+1)}(x) = \frac{\sigma_w}{\sqrt{d_l}} W^{(l+1)} \phi\left(z^{(l)}(x)\right) + \sigma_b b^{(l+1)}$$

## Notation

- Fully-connected L-layer network:

$$z^{(l+1)}(x) = \frac{\sigma_w}{\sqrt{d_l}} W^{(l+1)} \phi\left(z^{(l)}(x)\right) + \sigma_b b^{(l+1)}$$

- **Output:** $f_\theta(x) = z^{(L)}(x)$

# Notation

- Fully-connected L-layer network:

$$z^{(l+1)}(x) = \frac{\sigma_w}{\sqrt{d_l}} W^{(l+1)} \phi\left(z^{(l)}(x)\right) + \sigma_b b^{(l+1)}$$

- **Output:** $f_\theta(x) = z^{(L)}(x)$

- **Initialization:** $W_{ij}^{(l)} \sim \mathcal{N}(0,1)$ and $b_i^{(l)} \sim \mathcal{N}(0,1)$

# Notation

- Fully-connected L-layer network:

$$z^{(l+1)}(x) = \frac{\sigma_w}{\sqrt{d_l}} W^{(l+1)} \phi\left(z^{(l)}(x)\right) + \sigma_b b^{(l+1)}$$

- **Output:** $f_\theta(x) = z^{(L)}(x)$

- **Initialization:** $W_{ij}^{(l)} \sim \mathcal{N}(0,1)$ and $b_i^{(l)} \sim \mathcal{N}(0,1)$

- We can control **variances** $\sigma_w$, $\sigma_b$ and study how the network behaves when varying these

# Setup

## Setup

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

## Setup

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

- Stack features and targets into matrices:

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

# Recap: NNGP

# Recap: NNGP

- At **initialization:**

$$f_{\boldsymbol{\theta}(0)}(\cdot) \xrightarrow{d_1,\ldots,d_L \to \infty} \mathcal{GP}(0, \Sigma^{(L)})$$

# Recap: NNGP

- At **initialization:**

$$f_{\boldsymbol{\theta}(0)}(\cdot) \xrightarrow{d_1,\ldots,d_L \to \infty} \mathcal{GP}(0, \Sigma^{(L)})$$

where we have the recursion:

1) $\Sigma^{(1)}(\boldsymbol{x}, \boldsymbol{x}') = \boldsymbol{x}^T \boldsymbol{x}'$

$$\Sigma^{(\ell)}\big|_{x,x'} = \begin{pmatrix} \Sigma^{(\ell)}(x,x) & \Sigma^{(\ell)}(x,x') \\ \Sigma(x',x) & \Sigma^{(\ell)}(x',x') \end{pmatrix}$$

2) $\Sigma^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}') = \sigma_w^2 \mathbb{E}_{(\boldsymbol{z} \sim \mathcal{N}(\boldsymbol{0}, \Sigma^{(l)}|_{\boldsymbol{x},\boldsymbol{x}'}))} \Big[ \phi(z_1)\phi(z_2) \Big] + \sigma_b^2$

# Recap: NNGP

- At **initialization:**

$$f_{\boldsymbol{\theta}(0)}(\cdot) \xrightarrow{d_1,\dots,d_L \to \infty} \mathcal{GP}(0, \Sigma^{(L)})$$

  where we have the recursion:

  1) $\Sigma^{(1)}(\boldsymbol{x}, \boldsymbol{x}') = \boldsymbol{x}^T \boldsymbol{x}'$

  2) $\Sigma^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}') = \sigma_w^2 \mathbb{E}_{(\boldsymbol{z} \sim \mathcal{N}(\boldsymbol{0}, \Sigma^{(l)}|_{\boldsymbol{x}, \boldsymbol{x}'}))} \left[ \phi(z_1) \phi(z_2) \right] + \sigma_b^2$

$$\Sigma^{(l)}\big|_{\mathsf{x},\mathsf{x}'} = \begin{pmatrix} \Sigma^{(l)}(\mathsf{x},\mathsf{x}) & \Sigma^{(l)}(\mathsf{x},\mathsf{x}') \\ \Sigma^{(l)}(\mathsf{x}',\mathsf{x}) & \Sigma^{(l)}(\mathsf{x}',\mathsf{x}') \end{pmatrix}$$

- If network has **multiple** outputs:

$$f_{\boldsymbol{\theta},i} \perp f_{\boldsymbol{\theta},j}$$

# Recap: Neural Tangent Kernel (I)

# Recap: Neural Tangent Kernel (I)

- Empirical **NTK**:

$$\hat{\Theta}_t^{(L)}(\mathbf{x}, \mathbf{x}') = \left(\nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}(t)}(\mathbf{x})\right)^T \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}(t)}(\mathbf{x}')$$

# Recap: Neural Tangent Kernel (I)

- Empirical **NTK**:

$$\hat{\Theta}_t^{(L)}(\boldsymbol{x}, \boldsymbol{x}') = \left(\nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}(t)}(\boldsymbol{x})\right)^T \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}(t)}(\boldsymbol{x}')$$

- Governs the **dynamics**:

$$\partial_t f_{\boldsymbol{\theta}(t)}(\boldsymbol{x}) = -\hat{\Theta}_t^{(L)}(\boldsymbol{x}, \boldsymbol{X})^T \frac{\partial L(f_{\boldsymbol{\theta}(t)}(\boldsymbol{X}), \boldsymbol{y})}{\partial f_{\boldsymbol{\theta}(t)}(\boldsymbol{X})}$$

# Recap: Neural Tangent Kernel (I)

- Empirical **NTK**:

$$\hat{\Theta}_t^{(L)}(\boldsymbol{x}, \boldsymbol{x}') = \left(\nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}(t)}(\boldsymbol{x})\right)^T \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}(t)}(\boldsymbol{x}')$$

- Governs the **dynamics**:

$$\partial_t f_{\boldsymbol{\theta}(t)}(\boldsymbol{x}) = -\hat{\Theta}_t^{(L)}(\boldsymbol{x}, \boldsymbol{X})^T \frac{\partial L(f_{\boldsymbol{\theta}(t)}(\boldsymbol{X}), \boldsymbol{y})}{\partial f_{\boldsymbol{\theta}(t)}(\boldsymbol{X})}$$

- **Infinite-width** limit:

*Deterministic and time-independent*

$$\hat{\Theta}_t^{(L)}(\boldsymbol{x}, \boldsymbol{x}') \xrightarrow{d_1, \ldots, d_L \to \infty} \Theta^{(L)}(\boldsymbol{x}, \boldsymbol{x}')$$

# Recap: Neural Tangent Kernel (II)

# Recap: Neural Tangent Kernel (II)

- Again recursive structure:

  1) $\Theta^{(1)}(\boldsymbol{x}, \boldsymbol{x}') = \Sigma^{(1)}(\boldsymbol{x}, \boldsymbol{x}') = \boldsymbol{x}^T \boldsymbol{x}'$

  2) $\Theta^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}') = \Theta^{(l)}(\boldsymbol{x}, \boldsymbol{x}') \dot{\Sigma}^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}') + \Sigma^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}')$

$$\dot{\Sigma}^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}') = E_{\tilde{z} \sim \mathcal{N}\left(0, \Sigma^{(l+1)}\big|_{\boldsymbol{x}, \boldsymbol{x}'}\right)}\left[\dot{\phi}(\tilde{z}_1)\dot{\phi}(\tilde{z}_2)\right]$$

# Recap: Neural Tangent Kernel (II)

- Again recursive structure:

  1) $\Theta^{(1)}(\boldsymbol{x}, \boldsymbol{x}') = \Sigma^{(1)}(\boldsymbol{x}, \boldsymbol{x}') = \boldsymbol{x}^T \boldsymbol{x}'$

  2) $\Theta^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}') = \Theta^{(l)}(\boldsymbol{x}, \boldsymbol{x}')\dot{\Sigma}^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}') + \Sigma^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}')$

  $$\dot{\Sigma}^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}') = E_{\boldsymbol{z} \sim \mathcal{N}\left(0, \Sigma^{(l+1)}\big|_{\boldsymbol{x}, \boldsymbol{x}'}\right)}\left[\dot{\phi}(z_1)\dot{\phi}(z_2)\right]$$

- **Closed-form** expression for NN at $t = \infty$ and squared loss:

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

# Recap: Neural Tangent Kernel (II)

- Again recursive structure:

  1) $\Theta^{(1)}(\boldsymbol{x}, \boldsymbol{x}') = \Sigma^{(1)}(\boldsymbol{x}, \boldsymbol{x}') = \boldsymbol{x}^T \boldsymbol{x}'$

  2) $\Theta^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}') = \Theta^{(l)}(\boldsymbol{x}, \boldsymbol{x}')\dot{\Sigma}^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}') + \Sigma^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}')$

  $$\dot{\Sigma}^{(l+1)}(\underline{x}, \underline{x}') = \mathop{E}_{\tilde{z} \sim \mathcal{N}\left(0, \, \Sigma^{(l+1)}|_{\underline{x}, \underline{x}'}\right)}\left[\dot{\phi}(\tilde{z}_1)\dot{\phi}(\tilde{z}_2)\right]$$

- **Closed-form** expression for NN at $t = \infty$ and squared loss:

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

- We are assuming $f_0(\boldsymbol{x}) = 0$ here

## What Have We Considered So Far?

- We considered very **overparametrized** networks:

$$\#params >> \#samples$$

# What Have We Considered So Far?

- We considered very **overparametrized** networks:

$$\#params >> \#samples$$

- Overparametrization achieved with **huge widths**

$$d_1, \ldots, d_L \to \infty$$

# What Have We Considered So Far?

- We considered very **overparametrized** networks:

$$\#params >> \#samples$$

- Overparametrization achieved with **huge widths**

$$d_1, \ldots, d_L \to \infty$$

- What happens if we let **depth** $L$ go to infinity as well?

# Large-Depth Limit

## Large-Depth Limit

**Two** kinds of limits:

## Large-Depth Limit

**Two** kinds of limits:

- Take width and depth to infinity **simultaneously**

## Large-Depth Limit

**Two** kinds of limits:

- Take width and depth to infinity **simultaneously**

- Take width limit **first** and **then** consider depth limit

## Large-Depth Limit

**Two** kinds of limits:

- Take width and depth to infinity **simultaneously**

- Take width limit **first** and **then** consider depth limit

- Two **different** results!

# Simultaneous Limit

## Simultaneous Limit

- Very **difficult** to analyze!

## Simultaneous Limit

- Very **difficult** to analyze!

- "*Finite Depth and Width Corrections to the Neural Tangent Kernel*" analyzes this setting

# Simultaneous Limit

- Very **difficult** to analyze!

- "*Finite Depth and Width Corrections to the Neural Tangent Kernel*" analyzes this setting

- They are only able to analyze **diagonal** entries $\Theta^{(\infty)}(\boldsymbol{x}, \boldsymbol{x})$

## Simultaneous Limit

- Very **difficult** to analyze!

- "*Finite Depth and Width Corrections to the Neural Tangent Kernel*" analyzes this setting

- They are only able to analyze **diagonal** entries $\Theta^{(\infty)}(\boldsymbol{x}, \boldsymbol{x})$

- They prove that **SGD** induces change $\Delta\Theta^{(\infty)}(\boldsymbol{x}, \boldsymbol{x}')$ with

$$\text{var}(\Delta\Theta^{(\infty)}(\boldsymbol{x}, \boldsymbol{x}')) > 0$$

## Simultaneous Limit

- Very **difficult** to analyze!

- "*Finite Depth and Width Corrections to the Neural Tangent Kernel*" analyzes this setting

- They are only able to analyze **diagonal** entries $\Theta^{(\infty)}(x, x)$

- They prove that **SGD** induces change $\Delta\Theta^{(\infty)}(x, x')$ with

$$\text{var}(\Delta\Theta^{(\infty)}(x, x')) > 0$$

  $\implies$ **Feature-learning** happening!

# Sequential Limit

# Sequential Limit

- **Easier** to analyze!

# Sequential Limit

- **Easier** to analyze!

- Study the **dynamical systems**:

  1 $\Sigma^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}') = \sigma_w^2 \mathbb{E}_{\left(\boldsymbol{z} \sim \mathcal{N}(\boldsymbol{0}, \Sigma^{(l)}|_{\boldsymbol{x}, \boldsymbol{x}'})\right)}\left[\phi(z_1)\phi(z_2)\right] + \sigma_b^2$

  2 $\Theta^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}') = \Theta^{(l)}(\boldsymbol{x}, \boldsymbol{x}')\dot{\Sigma}^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}') + \Sigma^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}')$

# Sequential Limit

- **Easier** to analyze!

- Study the **dynamical systems**:

  1 $\Sigma^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}') = \sigma_w^2 \mathbb{E}_{\left(\boldsymbol{z} \sim \mathcal{N}(\boldsymbol{0}, \Sigma^{(l)}|_{\boldsymbol{x}, \boldsymbol{x}'})\right)} \left[\phi(z_1)\phi(z_2)\right] + \sigma_b^2$

  2 $\Theta^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}') = \Theta^{(l)}(\boldsymbol{x}, \boldsymbol{x}')\dot{\Sigma}^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}') + \Sigma^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}')$

- How do they behave as $l \to \infty$ as a function of $\sigma_w$ and $\sigma_b$?

# Sequential Limit

- **Easier** to analyze!

- Study the **dynamical systems**:

  1 $\Sigma^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}') = \sigma_w^2 \mathbb{E}_{(\boldsymbol{z} \sim \mathcal{N}(\boldsymbol{0}, \Sigma^{(l)}|_{\boldsymbol{x}, \boldsymbol{x}'}))} \left[ \phi(z_1)\phi(z_2) \right] + \sigma_b^2$

  2 $\Theta^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}') = \Theta^{(l)}(\boldsymbol{x}, \boldsymbol{x}')\dot{\Sigma}^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}') + \Sigma^{(l+1)}(\boldsymbol{x}, \boldsymbol{x}')$

- How do they behave as $l \to \infty$ as a function of $\sigma_w$ and $\sigma_b$?

- What happens to

$$P(\boldsymbol{\Theta}^{(L)}) = \left( \boldsymbol{\Theta}^{(L)}(\boldsymbol{x}, \boldsymbol{X}) \right)^T \left( \boldsymbol{\Theta}^{(L)}(\boldsymbol{X}, \boldsymbol{X}) \right)^{-1}$$

**This Work**

## This Work

- Consider the **sequential limit**

## This Work

- Consider the **sequential limit**

- **Incremental**:
  - "*Exponential expressivity in deep neural networks through transient chaos*"
  - "*Deep Information Propagation*"

# This Work

- Consider the **sequential limit**

- **Incremental**:
  - "*Exponential expressivity in deep neural networks through transient chaos*"
  - "*Deep Information Propagation*"

- Connect **chaotic** and **ordered** regimes with generalization and trainability

## This Work

- Consider the **sequential limit**

- **Incremental**:
  - "*Exponential expressivity in deep neural networks through transient chaos*"
  - "*Deep Information Propagation*"

- Connect **chaotic** and **ordered** regimes with generalization and trainability

- How to characterize **generalization** and **trainability**?

# Trainability (I)

# Trainability (I)

- Recall that we have closed-form training dynamics:

$$f_t(\boldsymbol{x}) = \boldsymbol{\Theta}^{(L)}(\boldsymbol{x}, \boldsymbol{X}) \left( \boldsymbol{\Theta}^{(L)}(\boldsymbol{X}, \boldsymbol{X}) \right)^{-1} \left( \mathbb{1} - e^{-t\eta\boldsymbol{\Theta}^{(L)}(\boldsymbol{X}, \boldsymbol{X})} \right) \boldsymbol{y}$$

## Trainability (I)

- Recall that we have closed-form training dynamics:

$$f_t(\boldsymbol{x}) = \boldsymbol{\Theta}^{(L)}(\boldsymbol{x}, \boldsymbol{X}) \left(\boldsymbol{\Theta}^{(L)}(\boldsymbol{X}, \boldsymbol{X})\right)^{-1} \left(\mathbb{1} - e^{-t\eta\boldsymbol{\Theta}^{(L)}(\boldsymbol{X},\boldsymbol{X})}\right) \boldsymbol{y}$$

- On the **training** set we have:

$$f_t(\boldsymbol{X}) = \left(\mathbb{1} - e^{-t\eta\boldsymbol{\Theta}^{(L)}(\boldsymbol{X},\boldsymbol{X})}\right) \boldsymbol{y}$$

# Trainability (I)

- Recall that we have closed-form training dynamics:

$$f_t(\boldsymbol{x}) = \Theta^{(L)}(\boldsymbol{x}, \boldsymbol{X}) \left( \Theta^{(L)}(\boldsymbol{X}, \boldsymbol{X}) \right)^{-1} \left( \mathbb{1} - e^{-t\eta\Theta^{(L)}(\boldsymbol{X}, \boldsymbol{X})} \right) \boldsymbol{y}$$

- On the **training** set we have:

$$f_t(\boldsymbol{X}) = \left( \mathbb{1} - e^{-t\eta\Theta^{(L)}(\boldsymbol{X}, \boldsymbol{X})} \right) \boldsymbol{y}$$

- **Diagonalize** NTK: $\Theta^{(L)}(\boldsymbol{X}, \boldsymbol{X}) = \boldsymbol{U}^T \boldsymbol{\Lambda} \boldsymbol{U}$ which leads to

$$\tilde{f}_{t,i}(\boldsymbol{X}) = \left( \mathbb{1} - e^{-t\eta\lambda_i} \right) \tilde{y}_i$$

rotated labels

rotated predictions

# Trainability (II)

# Trainability (II)

- **Maximum** feasible learning rate $\eta \sim \frac{2}{\lambda_1}$ (Recall Seyed's presentation)

# Trainability (II)

- **Maximum** feasible learning rate $\eta \sim \frac{2}{\lambda_1}$ (Recall Seyed's presentation)

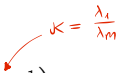- Plugging this into the dynamics:

$$\tilde{f}_{t,i}(\boldsymbol{X}) = \left( \mathbb{1} - e^{-2t\frac{\lambda_i}{\lambda_1}} \right) \tilde{y}_i$$

# Trainability (II)

- **Maximum** feasible learning rate $\eta \sim \frac{2}{\lambda_1}$ (Recall Seyed's presentation)

- Plugging this into the dynamics:

$$\tilde{f}_{t,i}(\boldsymbol{X}) = \left(1 - e^{-2t\frac{\lambda_i}{\lambda_1}}\right)\tilde{y}_i$$

- **Hardest** component to learn follows

$$\tilde{f}_{t,m}(\boldsymbol{X}) = \left(1 - e^{-2t\frac{\lambda_m}{\lambda_1}}\right)\tilde{y}_m = \left(1 - e^{-2t\kappa^{-1}}\right)\tilde{y}_m$$

$$\kappa = \frac{\lambda_1}{\lambda_m}$$

# Trainability (II)

- **Maximum** feasible learning rate $\eta \sim \frac{2}{\lambda_1}$ (Recall Seyed's presentation)

- Plugging this into the dynamics:

$$\tilde{f}_{t,i}(\boldsymbol{X}) = \left(\mathbb{1} - e^{-2t\frac{\lambda_i}{\lambda_1}}\right)\tilde{y}_i$$

- **Hardest** component to learn follows

$$\kappa = \frac{\lambda_1}{\lambda_m}$$

$$\tilde{f}_{t,m}(\boldsymbol{X}) = \left(\mathbb{1} - e^{-2t\frac{\lambda_m}{\lambda_1}}\right)\tilde{y}_m = \left(\mathbb{1} - e^{-2t\kappa^{-1}}\right)\tilde{y}_m$$

- **Hence:** If $\kappa \xrightarrow{L\to\infty} \infty$:

$$\tilde{f}_{t,m}^{(\infty)}(\boldsymbol{X}) = 0 \neq \tilde{y}_m \quad \forall t$$

# Trainability (II)

- **Maximum** feasible learning rate $\eta \sim \frac{2}{\lambda_1}$ (Recall Seyed's presentation)

- Plugging this into the dynamics:

$$\tilde{f}_{t,i}(\boldsymbol{X}) = \left( \mathbb{1} - e^{-2t\frac{\lambda_i}{\lambda_1}} \right) \tilde{y}_i$$

- **Hardest** component to learn follows

$$\tilde{f}_{t,m}(\boldsymbol{X}) = \left( \mathbb{1} - e^{-2t\frac{\lambda_m}{\lambda_1}} \right) \tilde{y}_m = \left( \mathbb{1} - e^{-2t\kappa^{-1}} \right) \tilde{y}_m$$

$$\kappa = \frac{\lambda_1}{\lambda_m}$$

- **Hence:** If $\kappa \xrightarrow{L \to \infty} \infty$:

$$\tilde{f}_{t,m}^{(\infty)}(\boldsymbol{X}) = 0 \neq \tilde{y}_m \quad \forall t$$

$\implies$ **Cannot** learn training set

# Generalization

## Generalization

- To study the **generalization** ability, look at $P(\Theta^{(\infty)})$

## Generalization

- To study the **generalization** ability, look at $P(\Theta^{(\infty)})$

- The limit will be a **function** of $(\sigma_w, \sigma_b)$

# Generalization

- To study the **generalization** ability, look at $P(\Theta^{(\infty)})$

- The limit will be a **function** of $(\sigma_w, \sigma_b)$

- Different regimes (chaotic and ordered) in the $(\sigma_w, \sigma_b)$ plane where $P(\Theta^{(\infty)})$ is **data-dependent** or **data-independent**

## Generalization

- To study the **generalization** ability, look at $P(\Theta^{(\infty)})$

- The limit will be a **function** of $(\sigma_w, \sigma_b)$

- Different regimes (chaotic and ordered) in the $(\sigma_w, \sigma_b)$ plane where $P(\Theta^{(\infty)})$ is **data-dependent** or **data-independent**

- Data-dependence of course **not** best measure to study generalization performance

# Goals

## Goals

- Characterize evolution of $\kappa^{(l)}$ to study trainability

  $\implies$ We hence need to control the spectrum of $\Theta^{(l)}$

## Goals

- Characterize evolution of $\kappa^{(l)}$ to study trainability

  $\implies$ We hence need to control the spectrum of $\Theta^{(l)}$

- Study trajectory of $P(\Theta^{(l)})$ to understand generalization

# NNGP and its Fixed Points

## NNGP and its Fixed Points

- **Assume :** Diagonal entries $q_{aa}^{(l)} := \Sigma^{(l)}(\boldsymbol{x}_a, \boldsymbol{x}_a)$ converge way quicker than off-diagonal terms $q_{ab}^{(l)} := \Sigma^{(l)}(\boldsymbol{x}_a, \boldsymbol{x}_b)$

# NNGP and its Fixed Points

- **Assume :** Diagonal entries $q_{aa}^{(l)} := \Sigma^{(l)}(x_a, x_a)$ converge way quicker than off-diagonal terms $q_{ab}^{(l)} := \Sigma^{(l)}(x_a, x_b)$

  **Approximate** $q_{aa}^{(l)} \approx q^* = \lim_{l \to \infty} q_{aa}^{(l)}$ in all calculations to follow.

# NNGP and its Fixed Points

- **Assume :** Diagonal entries $q_{aa}^{(l)} := \Sigma^{(l)}(\mathbf{x}_a, \mathbf{x}_a)$ converge way quicker than off-diagonal terms $q_{ab}^{(l)} := \Sigma^{(l)}(\mathbf{x}_a, \mathbf{x}_b)$

  **Approximate** $q_{aa}^{(l)} \approx q^* = \lim_{l \to \infty} q_{aa}^{(l)}$ in all calculations to follow.

  *Correlation between representations*

- Introduce $c_{ab}^{(l)} = \dfrac{q_{ab}^{(l)}}{\sqrt{q_{aa}^{(l)} q_{bb}^{(l)}}} = \dfrac{q_{ab}^{(l)}}{q^*}$ and write off-diagonal terms:

# NNGP and its Fixed Points

- **Assume :** Diagonal entries $q_{aa}^{(l)} := \Sigma^{(l)}(\boldsymbol{x}_a, \boldsymbol{x}_a)$ converge way quicker than off-diagonal terms $q_{ab}^{(l)} := \Sigma^{(l)}(\boldsymbol{x}_a, \boldsymbol{x}_b)$

  **Approximate** $q_{aa}^{(l)} \approx q^* = \lim_{l \to \infty} q_{aa}^{(l)}$ in all calculations to follow.

  *Correlation between representations*

- Introduce $c_{ab}^{(l)} = \dfrac{q_{ab}^{(l)}}{\sqrt{q_{aa}^{(l)} q_{bb}^{(l)}}} = \dfrac{q_{ab}^{(l)}}{q^*}$ and write off-diagonal terms:

$$q_{ab}^{(l+1)} = \sigma_w^2 \mathbb{E}_{\left(\boldsymbol{z} \sim \mathcal{N}\left(\boldsymbol{0}, \Sigma^{(l)}\big|_{\boldsymbol{x}, \boldsymbol{x}'}\right)\right)} \left[\phi(z_1)\phi(z_2)\right] + \sigma_b^2$$

# NNGP and its Fixed Points

- **Assume :** Diagonal entries $q_{aa}^{(l)} := \Sigma^{(l)}(x_a, x_a)$ converge way quicker than off-diagonal terms $q_{ab}^{(l)} := \Sigma^{(l)}(x_a, x_b)$

  **Approximate** $q_{aa}^{(l)} \approx q^* = \lim_{l \to \infty} q_{aa}^{(l)}$ in all calculations to follow.

  *Correlation between representations*

- Introduce $c_{ab}^{(l)} = \dfrac{q_{ab}^{(l)}}{\sqrt{q_{aa}^{(l)} q_{bb}^{(l)}}} = \dfrac{q_{ab}^{(l)}}{q^*}$ and write off-diagonal terms:

$$q_{ab}^{(l+1)} = \sigma_w^2 \mathbb{E}_{\left(z \sim \mathcal{N}\left(0, \Sigma^{(l)}\big|_{x, x'}\right)\right)} \left[\phi(z_1)\phi(z_2)\right] + \sigma_b^2$$

*Gaussian measure:* $\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z_1^2} \, dz_1$

$$= \sigma_w^2 \int \phi\left(\sqrt{q^*} z_1\right) \phi\left(\sqrt{q^*}(c_{ab}^{(l)} z_1 + \sqrt{1 - (c_{ab}^{(l)})^2} z_2)\right) \mathcal{D}z_1 \mathcal{D}z_2 + \sigma_b^2$$

# NNGP and its Fixed Points

- **Assume :** Diagonal entries $q_{aa}^{(l)} := \Sigma^{(l)}(\boldsymbol{x}_a, \boldsymbol{x}_a)$ converge way quicker than off-diagonal terms $q_{ab}^{(l)} := \Sigma^{(l)}(\boldsymbol{x}_a, \boldsymbol{x}_b)$

  **Approximate** $q_{aa}^{(l)} \approx q^* = \lim_{l \to \infty} q_{aa}^{(l)}$ in all calculations to follow.

  Correlation between representations

- Introduce $c_{ab}^{(l)} = \dfrac{q_{ab}^{(l)}}{\sqrt{q_{aa}^{(l)} q_{bb}^{(l)}}} = \dfrac{q_{ab}^{(l)}}{q^*}$ and write off-diagonal terms:

$$q_{ab}^{(l+1)} = \sigma_w^2 \mathbb{E}_{\left(\boldsymbol{z} \sim \mathcal{N}\left(\boldsymbol{0}, \Sigma^{(l)}\big|_{\boldsymbol{x}, \boldsymbol{x}'}\right)\right)} \left[\phi(z_1)\phi(z_2)\right] + \sigma_b^2$$

Gaussian measure: $\dfrac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z_1^2} dz_1$

$$= \sigma_w^2 \int \phi\left(\sqrt{q^*} z_1\right) \phi\left(\sqrt{q^*}(c_{ab}^{(l)} z_1 + \sqrt{1 - (c_{ab}^{(l)})^2} z_2)\right) \mathcal{D}z_1 \mathcal{D}z_2 + \sigma_b^2$$

$$= \sigma_w^2 \int \phi\left(u_1(z_1)\right) \phi\left(u_2(z_1, z_2, c_{ab}^{(l)})\right) \mathcal{D}z_1 \mathcal{D}z_2 + \sigma_b^2$$

# Recursive Correlation

# Recursive Correlation

- Define the function

$$f(c) = \frac{1}{q^*} \left( \sigma_w^2 \int \phi\left(u_1(z_1)\right) \phi\left(u_2(z_1, z_2, c)\right) \mathcal{D}z_1 \mathcal{D}z_2 + \sigma_b^2 \right)$$

# Recursive Correlation

- Define the function

$$f(c) = \frac{1}{q^*} \left( \sigma_w^2 \int \phi\left(u_1(z_1)\right) \phi\left(u_2(z_1, z_2, c)\right) \mathcal{D}z_1 \mathcal{D}z_2 + \sigma_b^2 \right)$$

- **Observe:** $f(1) = 1 \implies c^* = 1$ is a fixed point of $f$

## Recursive Correlation

- Define the function

$$f(c) = \frac{1}{q^*} \left( \sigma_w^2 \int \phi\left(u_1(z_1)\right) \phi\left(u_2(z_1, z_2, c)\right) \mathcal{D}z_1 \mathcal{D}z_2 + \sigma_b^2 \right)$$

- **Observe:** $f(1) = 1 \implies c^* = 1$ is a fixed point of $f$

- **Attracting** or **repulsive** fixed point?

# Recursive Correlation

- Define the function

$$f(c) = \frac{1}{q^*}\left(\sigma_w^2 \int \phi\left(u_1(z_1)\right)\phi\left(u_2(z_1, z_2, c)\right)\mathcal{D}z_1\mathcal{D}z_2 + \sigma_b^2\right)$$

- **Observe:** $f(1) = 1 \implies c^* = 1$ is a fixed point of $f$

- **Attracting** or **repulsive** fixed point?

- Calculus 101: Equivalent to $f'(1) < 1$ or $f'(1) > 1$?

## Recursive Correlation

- Define the function

$$f(c) = \frac{1}{q^*} \left( \sigma_w^2 \int \phi\left(u_1(z_1)\right) \phi\left(u_2(z_1, z_2, c)\right) \mathcal{D}z_1 \mathcal{D}z_2 + \sigma_b^2 \right)$$

- **Observe:** $f(1) = 1 \implies c^* = 1$ is a fixed point of $f$

- **Attracting** or **repulsive** fixed point?

- Calculus 101: Equivalent to $f'(1) < 1$ or $f'(1) > 1$?

- Turns out: $f'(1) = \sigma_w^2 \int \left[\phi'(\sqrt{q^*}z)\right]^2 \mathcal{D}z$

  $\implies$ Depends on $\sigma_w$ and $\sigma_b$

**What are Chaotic and Ordered Regimes?**

# What are Chaotic and Ordered Regimes?

- **Ordered**: $f'(1) < 1 \implies$ Correlation $c_{ab}^{(l)}$ between hidden representations of inputs $x_a$ and $x_b$ converges to $c^* = 1$:

$$\lim_{l \to \infty} c_{ab}^{(l)} = 1$$

## What are Chaotic and Ordered Regimes?

- **Ordered**: $f'(1) < 1 \implies$ Correlation $c_{ab}^{(l)}$ between hidden representations of inputs $x_a$ and $x_b$ converges to $c^* = 1$:

$$\lim_{l \to \infty} c_{ab}^{(l)} = 1$$

  This means that two outputs always **collapse onto** each other in the infinite depth limit

# What are Chaotic and Ordered Regimes?

- **Ordered**: $f'(1) < 1 \implies$ Correlation $c_{ab}^{(l)}$ between hidden representations of inputs $x_a$ and $x_b$ converges to $c^* = 1$:

$$\lim_{l \to \infty} c_{ab}^{(l)} = 1$$

  This means that two outputs always **collapse onto** each other in the infinite depth limit

- **Chaotic**: $f'(1) > 1 \implies c^* = 1$ is repulsive

## What are Chaotic and Ordered Regimes?

- **Ordered**: $f'(1) < 1 \implies$ Correlation $c_{ab}^{(l)}$ between hidden representations of inputs $x_a$ and $x_b$ converges to $c^* = 1$:

$$\lim_{l \to \infty} c_{ab}^{(l)} = 1$$

  This means that two outputs always **collapse onto** each other in the infinite depth limit

- **Chaotic**: $f'(1) > 1 \implies c^* = 1$ is repulsive

  Since $c_{ab}^{(l)} < 1$:

$$\lim_{l \to \infty} c_{ab}^{(l)} < 1$$

## What are Chaotic and Ordered Regimes?

- **Ordered**: $f'(1) < 1 \implies$ Correlation $c_{ab}^{(l)}$ between hidden representations of inputs $x_a$ and $x_b$ converges to $c^* = 1$:

$$\lim_{l \to \infty} c_{ab}^{(l)} = 1$$

  This means that two outputs always **collapse onto** each other in the infinite depth limit

- **Chaotic**: $f'(1) > 1 \implies c^* = 1$ is repulsive

  Since $c_{ab}^{(l)} < 1$:

$$\lim_{l \to \infty} c_{ab}^{(l)} < 1$$

  Two outputs thus become more and more **dissimilar** to each other

# Restate the Dynamics

# Restate the Dynamics

- Dynamics in terms of new notation:

# Restate the Dynamics

- Dynamics in terms of new notation:

**1a)** $q_{ab}^{(l+1)} = \sigma_w^2 \mathcal{T}\left(q_{ab}^{(l)}\right) + \sigma_b^2$

**1b)** $q_{aa}^{(l+1)} = q^*$

$$= \mathbb{E}_{z \sim \mathcal{N}\left(0, \begin{pmatrix} q^* & q_{ab}^{(l)} \\ q_{ab}^{(l)} & q^* \end{pmatrix}\right)} \left[\phi(z_1)\phi(z_2)\right]$$

# Restate the Dynamics

- Dynamics in terms of new notation:

**1a)** $q_{ab}^{(l+1)} = \sigma_w^2 \mathcal{T}\left(q_{ab}^{(l)}\right) + \sigma_b^2$

**1b)** $q_{aa}^{(l+1)} = q^*$ $\qquad = \mathbb{E}_{z \sim \mathcal{N}\left(0, \begin{pmatrix} q^* & q_{ab}^{(l)} \\ q_{ab}^{(l)} & q^* \end{pmatrix}\right)}\left[\phi(z_1)\phi(z_2)\right]$

**2a)** $p_{ab}^{(l+1)} = q_{ab}^{(l+1)} + \sigma_w^2 \dot{\mathcal{T}}(q_{ab}^{(l)}) p_{ab}^{(l)}$

**2b)** $p_{aa}^{(l+1)} = q^* + \sigma_w^2 \dot{\mathcal{T}}(q^*) p_{aa}^{(l)} = q^* + f'(1) p_{aa}^{(l)}$

# Restate the Dynamics

- Dynamics in terms of new notation:

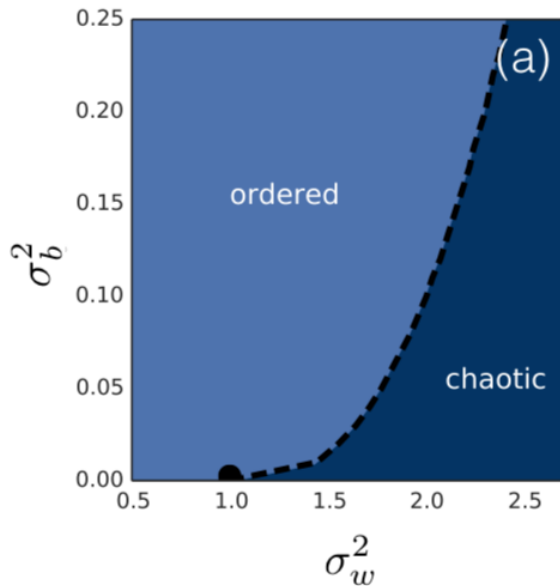**1a)** $q_{ab}^{(l+1)} = \sigma_w^2 \mathcal{T}\left(q_{ab}^{(l)}\right) + \sigma_b^2$

**1b)** $q_{aa}^{(l+1)} = q^*$
$$= \mathbb{E}_{z \sim \mathcal{N}\left(0, \begin{pmatrix} q^* & q_{ab}^{(l)} \\ q_{ab}^{(l)} & q^* \end{pmatrix}\right)}\left[\phi(z_1)\phi(z_2)\right]$$

**2a)** $p_{ab}^{(l+1)} = q_{ab}^{(l+1)} + \sigma_w^2 \dot{\mathcal{T}}(q_{ab}^{(l)})p_{ab}^{(l)}$

**2b)** $p_{aa}^{(l+1)} = q^* + \sigma_w^2 \dot{\mathcal{T}}(q^*)p_{aa}^{(l)} = q^* + f'(1)p_{aa}^{(l)}$

- Analyze the limits in the two different regimes

## An Example

# Chaotic Regime: Limits

# Chaotic Regime: Limits

**1.b)** $q_{aa}^{(l)} = q^*$ already given

# Chaotic Regime: Limits

**1.b)** $q_{aa}^{(l)} = q^*$ already given

**1.a) Recall**: $c_{ab}^{(l)} = \frac{q_{ab}^{(l)}}{q^*} \xrightarrow{l \to \infty} c^* < 1 \implies q_{ab}^{(l)} \xrightarrow{l \to \infty} c^* q^*$

# Chaotic Regime: Limits

**1.b)** $q_{aa}^{(l)} = q^*$ already given

**1.a) Recall**: $c_{ab}^{(l)} = \frac{q_{ab}^{(l)}}{q^*} \xrightarrow{l \to \infty} c^* < 1 \implies q_{ab}^{(l)} \xrightarrow{l \to \infty} c^* q^*$

**2.b)** $p_{aa}^{(l+1)} = q^* + f'(1) p_{aa}^{(l)}$

# Chaotic Regime: Limits

**1.b)** $q_{aa}^{(l)} = q^*$ already given

**1.a) Recall:** $c_{ab}^{(l)} = \frac{q_{ab}^{(l)}}{q^*} \xrightarrow{l \to \infty} c^* < 1 \implies q_{ab}^{(l)} \xrightarrow{l \to \infty} c^* q^*$

**2.b)** $p_{aa}^{(l+1)} = q^* + f'(1)p_{aa}^{(l)} = \cdots = \sum_{i=0}^{l+1} q^* f'(1)^i$

# Chaotic Regime: Limits

**1.b)** $q_{aa}^{(l)} = q^*$ already given

**1.a) Recall**: $c_{ab}^{(l)} = \frac{q_{ab}^{(l)}}{q^*} \xrightarrow{l \to \infty} c^* < 1 \implies q_{ab}^{(l)} \xrightarrow{l \to \infty} c^* q^*$

**2.b)** $p_{aa}^{(l+1)} = q^* + f'(1)p_{aa}^{(l)} = \cdots = \sum_{i=0}^{l+1} q^* f'(1)^i = \frac{f'(1)^{l+1} - 1}{f'(1) - 1} q^*$

# Chaotic Regime: Limits

**1.b)** $q_{aa}^{(l)} = q^*$ already given

**1.a) Recall**: $c_{ab}^{(l)} = \frac{q_{ab}^{(l)}}{q^*} \xrightarrow{l \to \infty} c^* < 1 \implies q_{ab}^{(l)} \xrightarrow{l \to \infty} c^* q^*$

**2.b)** $p_{aa}^{(l+1)} = q^* + f'(1)p_{aa}^{(l)} = \cdots = \sum_{i=0}^{l+1} q^* f'(1)^i = \frac{f'(1)^{l+1}-1}{f'(1)-1} q^*$

**Hence:** $p_{aa}^{(l+1)} = \frac{f'(1)^{l+1}-1}{f'(1)-1} q^*$

# Chaotic Regime: Limits

**1.b)** $q_{aa}^{(l)} = q^*$ already given

**1.a) Recall:** $c_{ab}^{(l)} = \frac{q_{ab}^{(l)}}{q^*} \xrightarrow{l \to \infty} c^* < 1 \implies q_{ab}^{(l)} \xrightarrow{l \to \infty} c^* q^*$

**2.b)** $p_{aa}^{(l+1)} = q^* + f'(1)p_{aa}^{(l)} = \cdots = \sum_{i=0}^{l+1} q^* f'(1)^i = \frac{f'(1)^{l+1} - 1}{f'(1) - 1} q^*$

**Hence:** $p_{aa}^{(l+1)} = \frac{f'(1)^{l+1} - 1}{f'(1) - 1} q^* \approx f'(1)^l q^* \xrightarrow{l \to \infty} \infty$

# Chaotic Regime: Limits

**1.b)** $q_{aa}^{(l)} = q^*$ already given

**1.a) Recall**: $c_{ab}^{(l)} = \frac{q_{ab}^{(l)}}{q^*} \xrightarrow{l \to \infty} c^* < 1 \implies q_{ab}^{(l)} \xrightarrow{l \to \infty} c^* q^*$

**2.b)** $p_{aa}^{(l+1)} = q^* + f'(1)p_{aa}^{(l)} = \cdots = \sum_{i=0}^{l+1} q^* f'(1)^i = \frac{f'(1)^{l+1}-1}{f'(1)-1}q^*$

**Hence:** $p_{aa}^{(l+1)} = \frac{f'(1)^{l+1}-1}{f'(1)-1}q^* \approx f'(1)^l q^* \xrightarrow{l \to \infty} \infty$

**2.a)** Taking the limit on both sides of **2a)**:

# Chaotic Regime: Limits

**1.b)** $q_{aa}^{(l)} = q^*$ already given

**1.a) Recall**: $c_{ab}^{(l)} = \frac{q_{ab}^{(l)}}{q^*} \xrightarrow{l \to \infty} c^* < 1 \implies q_{ab}^{(l)} \xrightarrow{l \to \infty} c^* q^*$

**2.b)** $p_{aa}^{(l+1)} = q^* + f'(1)p_{aa}^{(l)} = \cdots = \sum_{i=0}^{l+1} q^* f'(1)^i = \frac{f'(1)^{l+1} - 1}{f'(1) - 1} q^*$

**Hence:** $p_{aa}^{(l+1)} = \frac{f'(1)^{l+1} - 1}{f'(1) - 1} q^* \approx f'(1)^l q^* \xrightarrow{l \to \infty} \infty$

**2.a)** Taking the limit on both sides of **2a)**:

$$p_{ab}^* = q_{ab}^* + \sigma_w^2 \dot{\mathcal{T}}(q_{ab}^*) p_{ab}^*$$

# Chaotic Regime: Limits

**1.b)** $q_{aa}^{(l)} = q^*$ already given

**1.a) Recall**: $c_{ab}^{(l)} = \frac{q_{ab}^{(l)}}{q^*} \xrightarrow{l \to \infty} c^* < 1 \implies q_{ab}^{(l)} \xrightarrow{l \to \infty} c^* q^*$

**2.b)** $p_{aa}^{(l+1)} = q^* + f'(1)p_{aa}^{(l)} = \cdots = \sum_{i=0}^{l+1} q^* f'(1)^i = \frac{f'(1)^{l+1}-1}{f'(1)-1} q^*$

**Hence:** $p_{aa}^{(l+1)} = \frac{f'(1)^{l+1}-1}{f'(1)-1} q^* \approx f'(1)^l q^* \xrightarrow{l \to \infty} \infty$

**2.a)** Taking the limit on both sides of **2a)**:

$$p_{ab}^* = q_{ab}^* + \sigma_w^2 \dot{\mathcal{T}}(q_{ab}^*) p_{ab}^* \iff p_{ab}^* = \frac{q_{ab}^*}{1 - f'(c^*)}$$

# Chaotic Regime: Spectrum NNGP

## Chaotic Regime: Spectrum NNGP

- We can summarize as follows:

## Chaotic Regime: Spectrum NNGP

- We can summarize as follows:

$$\Sigma^{(\infty)}(\boldsymbol{X}, \boldsymbol{X}) = q^*(c^*\mathbf{1}\mathbf{1}^T + (1-c^*)\mathbb{1}_{n \times n}) = \begin{pmatrix} q^* & c^*q^* & \dots & c^*q^* \\ c^*q^* & q^* & \dots & c^*q^* \\ \vdots & & & \\ c^*q^* & \dots & q^* & c^*q^* \\ c^*q^* & \dots & c^*q^* & q^* \end{pmatrix}$$

## Chaotic Regime: Spectrum NNGP

- We can summarize as follows:

$$\Sigma^{(\infty)}(\boldsymbol{X}, \boldsymbol{X}) = q^*(c^*\boldsymbol{1}\boldsymbol{1}^T + (1-c^*)\mathbb{1}_{n\times n}) = \begin{pmatrix} q^* & c^*q^* & \dots & c^*q^* \\ c^*q^* & q^* & \dots & c^*q^* \\ \vdots & & & \\ c^*q^* & \dots & q^* & c^*q^* \\ c^*q^* & \dots & c^*q^* & q^* \end{pmatrix}$$

- Has eigenvector $\boldsymbol{v}_1 = (1, \dots, 1)$ with $\lambda_0 = q^*\left((n-1)c^* + 1\right)$

## Chaotic Regime: Spectrum NNGP

- We can summarize as follows:

$$\Sigma^{(\infty)}(\boldsymbol{X}, \boldsymbol{X}) = q^*(c^* \mathbf{1}\mathbf{1}^T + (1-c^*)\mathbb{1}_{n \times n}) = \begin{pmatrix} q^* & c^*q^* & \dots & c^*q^* \\ c^*q^* & q^* & \dots & c^*q^* \\ \vdots & & & \\ c^*q^* & \dots & q^* & c^*q^* \\ c^*q^* & \dots & c^*q^* & q^* \end{pmatrix}$$

- Has eigenvector $\boldsymbol{v}_1 = (1, \dots, 1)$ with $\lambda_0 = q^*\left((n-1)c^* + 1\right)$

- Others are $\boldsymbol{v}_i = (\boldsymbol{e}_i - \boldsymbol{e}_{i+1})$ and $\lambda_i = \lambda_{\text{bulk}} = q^*(1 - c^*)$

## Chaotic Regime: Spectrum NNGP

- We can summarize as follows:

$$\Sigma^{(\infty)}(\boldsymbol{X}, \boldsymbol{X}) = q^*(c^* \mathbf{11}^T + (1-c^*)\mathbb{1}_{n \times n}) = \begin{pmatrix} q^* & c^*q^* & \ldots & c^*q^* \\ c^*q^* & q^* & \ldots & c^*q^* \\ \vdots & & & \\ c^*q^* & \ldots & q^* & c^*q^* \\ c^*q^* & \ldots & c^*q^* & q^* \end{pmatrix}$$

- Has eigenvector $\boldsymbol{v}_1 = (1, \ldots, 1)$ with $\lambda_0 = q^*\left((n-1)c^* + 1\right)$

- Others are $\boldsymbol{v}_i = (\boldsymbol{e}_i - \boldsymbol{e}_{i+1})$ and $\lambda_i = \lambda_{\text{bulk}} = q^*(1-c^*)$

- $\kappa(\Sigma^{(\infty)}(\boldsymbol{X}, \boldsymbol{X})) = \frac{\lambda_0}{\lambda_{\text{bulk}}} = \frac{(n-1)c^*+1}{1-c^*}$

# Chaotic Regime: Spectrum NTK

# Chaotic Regime: Spectrum NTK

- $\frac{1}{p_{aa}^*} \Theta^{(\infty)}(\boldsymbol{X}, \boldsymbol{X}) = \mathbb{1}_{n \times n}$

# Chaotic Regime: Spectrum NTK

- $\frac{1}{p_{aa}^*} \Theta^{(\infty)}(\boldsymbol{X}, \boldsymbol{X}) = \mathbb{1}_{n \times n}$

- $\lambda_i \left( \Theta^{(\infty)}(\boldsymbol{X}, \boldsymbol{X}) \right) = \infty$

# Chaotic Regime: Spectrum NTK

- $\frac{1}{p_{aa}^*}\Theta^{(\infty)}(\boldsymbol{X}, \boldsymbol{X}) = \mathbb{1}_{n \times n}$

- $\lambda_i\left(\Theta^{(\infty)}(\boldsymbol{X}, \boldsymbol{X})\right) = \infty$

- $\kappa\left(\Theta^{(\infty)}(\boldsymbol{X}, \boldsymbol{X})\right) = 1 \implies$ NTK **remains trainable** for infinite-depth in the chaotic regime

# Ordered Regime: Limits

# Ordered Regime: Limits

- $q_{aa}^{(I)} = q^*$ already given

# Ordered Regime: Limits

- $q_{aa}^{(I)} = q^*$ already given

- **Recall**: $c_{ab}^{(I)} = \frac{q_{ab}^{(I)}}{q^*} \xrightarrow{I \to \infty} c^* = 1 \implies q_{ab}^{(I)} \xrightarrow{I \to \infty} q^*$

# Ordered Regime: Limits

- $q_{aa}^{(l)} = q^*$ already given

- **Recall**: $c_{ab}^{(l)} = \frac{q_{ab}^{(l)}}{q^*} \xrightarrow{l \to \infty} c^* = 1 \implies q_{ab}^{(l)} \xrightarrow{l \to \infty} q^*$

- Using $f'(1) < 1$, we can apply the geometric series

# Ordered Regime: Limits

- $q_{aa}^{(l)} = q^*$ already given

- **Recall**: $c_{ab}^{(l)} = \frac{q_{ab}^{(l)}}{q^*} \xrightarrow{l \to \infty} c^* = 1 \implies q_{ab}^{(l)} \xrightarrow{l \to \infty} q^*$

- Using $f'(1) < 1$, we can apply the geometric series

$$p_{aa}^{(l+1)} = q^* + f'(1)p_{aa}^{(l)}$$

## Ordered Regime: Limits

- $q_{aa}^{(l)} = q^*$ already given

- **Recall**: $c_{ab}^{(l)} = \frac{q_{ab}^{(l)}}{q^*} \xrightarrow{l \to \infty} c^* = 1 \implies q_{ab}^{(l)} \xrightarrow{l \to \infty} q^*$

- Using $f'(1) < 1$, we can apply the geometric series

$$p_{aa}^{(l+1)} = q^* + f'(1)p_{aa}^{(l)} = \cdots = \sum_{i=0}^{l+1} q^* f'(1)^i$$

## Ordered Regime: Limits

- $q_{aa}^{(l)} = q^*$ already given

- **Recall**: $c_{ab}^{(l)} = \frac{q_{ab}^{(l)}}{q^*} \xrightarrow{l \to \infty} c^* = 1 \implies q_{ab}^{(l)} \xrightarrow{l \to \infty} q^*$

- Using $f'(1) < 1$, we can apply the geometric series

$$p_{aa}^{(l+1)} = q^* + f'(1)p_{aa}^{(l)} = \cdots = \sum_{i=0}^{l+1} q^* f'(1)^i \xrightarrow{l \to \infty} \sum_{i=0}^{\infty} q^* \left(f'(1)\right)^i$$

# Ordered Regime: Limits

- $q_{aa}^{(l)} = q^*$ already given

- **Recall**: $c_{ab}^{(l)} = \frac{q_{ab}^{(l)}}{q^*} \xrightarrow{l \to \infty} c^* = 1 \implies q_{ab}^{(l)} \xrightarrow{l \to \infty} q^*$

- Using $f'(1) < 1$, we can apply the geometric series

$$p_{aa}^{(l+1)} = q^* + f'(1)p_{aa}^{(l)} = \cdots = \sum_{i=0}^{l+1} q^* f'(1)^i \xrightarrow{l \to \infty} \sum_{i=0}^{\infty} q^* \left(f'(1)\right)^i$$

$$= \frac{q^*}{1 - f'(1)}$$

## Ordered Regime: Limits

- $q_{aa}^{(l)} = q^*$ already given

- **Recall**: $c_{ab}^{(l)} = \frac{q_{ab}^{(l)}}{q^*} \xrightarrow{l \to \infty} c^* = 1 \implies q_{ab}^{(l)} \xrightarrow{l \to \infty} q^*$

- Using $f'(1) < 1$, we can apply the geometric series

$$p_{aa}^{(l+1)} = q^* + f'(1)p_{aa}^{(l)} = \cdots = \sum_{i=0}^{l+1} q^* f'(1)^i \xrightarrow{l \to \infty} \sum_{i=0}^{\infty} q^* \left(f'(1)\right)^i$$

$$= \frac{q^*}{1 - f'(1)}$$

- Taking the limit on both sides of **2a)**:

## Ordered Regime: Limits

- $q_{aa}^{(l)} = q^*$ already given

- **Recall**: $c_{ab}^{(l)} = \frac{q_{ab}^{(l)}}{q^*} \xrightarrow{l \to \infty} c^* = 1 \implies q_{ab}^{(l)} \xrightarrow{l \to \infty} q^*$

- Using $f'(1) < 1$, we can apply the geometric series

$$p_{aa}^{(l+1)} = q^* + f'(1)p_{aa}^{(l)} = \cdots = \sum_{i=0}^{l+1} q^* f'(1)^i \xrightarrow{l \to \infty} \sum_{i=0}^{\infty} q^* \left(f'(1)\right)^i$$

$$= \frac{q^*}{1 - f'(1)}$$

- Taking the limit on both sides of **2a)**:

$$p_{ab}^* = q_{ab}^* + \sigma_w^2 \dot{\mathcal{T}}(q^*)p_{ab}^*$$

## Ordered Regime: Limits

- $q_{aa}^{(l)} = q^*$ already given

- **Recall**: $c_{ab}^{(l)} = \frac{q_{ab}^{(l)}}{q^*} \xrightarrow{l \to \infty} c^* = 1 \implies q_{ab}^{(l)} \xrightarrow{l \to \infty} q^*$

- Using $f'(1) < 1$, we can apply the geometric series

$$p_{aa}^{(l+1)} = q^* + f'(1)p_{aa}^{(l)} = \cdots = \sum_{i=0}^{l+1} q^* f'(1)^i \xrightarrow{l \to \infty} \sum_{i=0}^{\infty} q^* \left(f'(1)\right)^i$$

$$= \frac{q^*}{1 - f'(1)}$$

- Taking the limit on both sides of **2a)**:

$$p_{ab}^* = q_{ab}^* + \sigma_w^2 \dot{\mathcal{T}}(q^*)p_{ab}^* \iff p_{ab}^* = \frac{q^*}{1 - f'(1)} = p^*$$

# Ordered Regime: Spectrum NNGP

## Ordered Regime: Spectrum NNGP

- We can summarize the NNGP as

$$\mathbf{\Sigma}^{(\infty)}(\boldsymbol{X}, \boldsymbol{X}) = q^* \mathbf{1}\mathbf{1}^T$$

## Ordered Regime: Spectrum NNGP

- We can summarize the NNGP as

$$\mathbf{\Sigma}^{(\infty)}(\mathbf{X}, \mathbf{X}) = q^* \mathbf{1}\mathbf{1}^T$$

- Again, $\mathbf{v}_1 = \mathbf{1}$ with $\lambda_1(\mathbf{\Sigma}^{(\infty)}(\mathbf{X}, \mathbf{X})) = nq^*$

## Ordered Regime: Spectrum NNGP

- We can summarize the NNGP as

$$\Sigma^{(\infty)}(\boldsymbol{X}, \boldsymbol{X}) = q^* \mathbf{1}\mathbf{1}^T$$

- Again, $\boldsymbol{v}_1 = \mathbf{1}$ with $\lambda_1(\Sigma^{(\infty)}(\boldsymbol{X}, \boldsymbol{X})) = nq^*$

- Moreover, $\boldsymbol{v}_i = \boldsymbol{e}_i - \boldsymbol{e}_{i+1}$ with $\lambda_{\mathsf{bulk}}(\Sigma^{(\infty)}(\boldsymbol{X}, \boldsymbol{X}) = 0$

# Ordered Regime: Spectrum NNGP

- We can summarize the NNGP as

$$\Sigma^{(\infty)}(\boldsymbol{X}, \boldsymbol{X}) = q^* \mathbf{1}\mathbf{1}^T$$

- Again, $\boldsymbol{v}_1 = \mathbf{1}$ with $\lambda_1(\Sigma^{(\infty)}(\boldsymbol{X}, \boldsymbol{X})) = nq^*$

- Moreover, $\boldsymbol{v}_i = \boldsymbol{e}_i - \boldsymbol{e}_{i+1}$ with $\lambda_{\text{bulk}}(\Sigma^{(\infty)}(\boldsymbol{X}, \boldsymbol{X}) = 0$

- We can calculate the condition number as

$$\kappa\left(\Sigma^{(\infty)}(\boldsymbol{X}, \boldsymbol{X})\right) = \frac{\lambda_1\left(\Sigma^{(\infty)}(\boldsymbol{X}, \boldsymbol{X})\right)}{\lambda_{\text{bulk}}\left(\Sigma^{(\infty)}(\boldsymbol{X}, \boldsymbol{X})\right)} = \infty$$

# Ordered Regime: Spectrum NTK

## Ordered Regime: Spectrum NTK

- We can express the infinite-depth NTK compactly via

$$\Theta^{(\infty)}(\boldsymbol{X}, \boldsymbol{X}) = p^* \mathbf{1}\mathbf{1}^T$$

# Ordered Regime: Spectrum NTK

- We can express the infinite-depth NTK compactly via

$$\mathbf{\Theta}^{(\infty)}(\boldsymbol{X}, \boldsymbol{X}) = p^* \mathbf{1}\mathbf{1}^T$$

- We get the same eigen-structure as for the NNGP

## Ordered Regime: Spectrum NTK

- We can express the infinite-depth NTK compactly via

$$\Theta^{(\infty)}(\boldsymbol{X}, \boldsymbol{X}) = p^* \mathbf{1}\mathbf{1}^T$$

- We get the same eigen-structure as for the NNGP

- $\lambda_1\left(\Theta^{(\infty)}\right) = np^*$

# Ordered Regime: Spectrum NTK

- We can express the infinite-depth NTK compactly via

$$\mathbf{\Theta}^{(\infty)}(\boldsymbol{X}, \boldsymbol{X}) = p^* \mathbf{1}\mathbf{1}^T$$

- We get the same eigen-structure as for the NNGP

- $\lambda_1\left(\mathbf{\Theta}^{(\infty)}\right) = np^*$

- $\lambda_{\text{bulk}}\left(\mathbf{\Theta}^{(\infty)}\right) = 0$

# Ordered Regime: Spectrum NTK

- We can express the infinite-depth NTK compactly via

$$\Theta^{(\infty)}(\boldsymbol{X}, \boldsymbol{X}) = p^* \mathbf{1}\mathbf{1}^T$$

- We get the same eigen-structure as for the NNGP

- $\lambda_1\left(\Theta^{(\infty)}\right) = np^*$

- $\lambda_{\text{bulk}}\left(\Theta^{(\infty)}\right) = 0$

- $\kappa\left(\Theta^{(\infty)}(\boldsymbol{X}, \boldsymbol{X})\right) = \infty$

# Ordered Regime: Spectrum NTK

- We can express the infinite-depth NTK compactly via

$$\mathbf{\Theta}^{(\infty)}(\mathbf{X}, \mathbf{X}) = p^* \mathbf{1}\mathbf{1}^T$$

- We get the same eigen-structure as for the NNGP

- $\lambda_1\left(\mathbf{\Theta}^{(\infty)}\right) = np^*$

- $\lambda_{\text{bulk}}\left(\mathbf{\Theta}^{(\infty)}\right) = 0$

- $\kappa\left(\mathbf{\Theta}^{(\infty)}(\mathbf{X}, \mathbf{X})\right) = \infty$

    $\implies$ Networks become **untrainable**!

# Finite Depth Correction

# Finite Depth Correction

- Define $\epsilon_{ab}^{(l)} = q_{ab}^{(l)} - q_{ab}^*$ and $\delta_{ab}^{(l)} = p_{ab}^{(l)} - p_{ab}^*$

# Finite Depth Correction

- Define $\epsilon_{ab}^{(l)} = q_{ab}^{(l)} - q_{ab}^*$ and $\delta_{ab}^{(l)} = p_{ab}^{(l)} - p_{ab}^*$

- Expanding and Taylor-approximating:

# Finite Depth Correction

- Define $\epsilon_{ab}^{(l)} = q_{ab}^{(l)} - q_{ab}^*$ and $\delta_{ab}^{(l)} = p_{ab}^{(l)} - p_{ab}^*$

- Expanding and Taylor-approximating:

$$q_{ab}^{(l+1)} = q_{ab}^* + \epsilon_{ab}^{(l+1)}$$

# Finite Depth Correction

- Define $\epsilon_{ab}^{(l)} = q_{ab}^{(l)} - q_{ab}^*$ and $\delta_{ab}^{(l)} = p_{ab}^{(l)} - p_{ab}^*$

- Expanding and Taylor-approximating:

$$q_{ab}^{(l+1)} = q_{ab}^* + \epsilon_{ab}^{(l+1)} = \sigma_w^2 \mathcal{T}\left(q_{ab}^{(l)}\right) + \sigma_b^2$$

# Finite Depth Correction

- Define $\epsilon_{ab}^{(l)} = q_{ab}^{(l)} - q_{ab}^*$ and $\delta_{ab}^{(l)} = p_{ab}^{(l)} - p_{ab}^*$

- Expanding and Taylor-approximating:

$$q_{ab}^{(l+1)} = q_{ab}^* + \epsilon_{ab}^{(l+1)} = \sigma_w^2 \mathcal{T}\left(q_{ab}^{(l)}\right) + \sigma_b^2 = \sigma_w^2 \mathcal{T}\left(q_{ab}^* + \epsilon_{ab}^{(l)}\right) + \sigma_b^2$$

# Finite Depth Correction

- Define $\epsilon_{ab}^{(l)} = q_{ab}^{(l)} - q_{ab}^*$ and $\delta_{ab}^{(l)} = p_{ab}^{(l)} - p_{ab}^*$

- Expanding and Taylor-approximating:

$$q_{ab}^{(l+1)} = q_{ab}^* + \epsilon_{ab}^{(l+1)} = \sigma_w^2 \mathcal{T}\left(q_{ab}^{(l)}\right) + \sigma_b^2 = \sigma_w^2 \mathcal{T}\left(q_{ab}^* + \epsilon_{ab}^{(l)}\right) + \sigma_b^2$$

$$\stackrel{Taylor}{=} \sigma_w^2 \mathcal{T}(q_{ab}^*) + \sigma_b^2 + \sigma_w^2 \dot{\mathcal{T}}(q_{ab}^*)\,\epsilon_{ab}^{(l)} + \mathcal{O}\left(\left(\epsilon_{ab}^{(l)}\right)^2\right)$$

# Finite Depth Correction

- Define $\epsilon_{ab}^{(l)} = q_{ab}^{(l)} - q_{ab}^*$ and $\delta_{ab}^{(l)} = p_{ab}^{(l)} - p_{ab}^*$

- Expanding and Taylor-approximating:

$$q_{ab}^{(l+1)} = q_{ab}^* + \epsilon_{ab}^{(l+1)} = \sigma_w^2 \mathcal{T}\left(q_{ab}^{(l)}\right) + \sigma_b^2 = \sigma_w^2 \mathcal{T}\left(q_{ab}^* + \epsilon_{ab}^{(l)}\right) + \sigma_b^2$$

$$\overset{\text{Taylor}}{=} \sigma_w^2 \mathcal{T}\left(q_{ab}^*\right) + \sigma_b^2 + \sigma_w^2 \dot{\mathcal{T}}\left(q_{ab}^*\right) \epsilon_{ab}^{(l)} + \mathcal{O}\left(\left(\epsilon_{ab}^{(l)}\right)^2\right)$$

$$= q_{ab}^* + \sigma_w^2 \dot{\mathcal{T}}\left(q_{ab}^*\right) \epsilon_{ab}^{(l)} + \mathcal{O}\left(\left(\epsilon_{ab}^{(l)}\right)^2\right)$$

# Finite Depth Correction

- Define $\epsilon_{ab}^{(l)} = q_{ab}^{(l)} - q_{ab}^*$ and $\delta_{ab}^{(l)} = p_{ab}^{(l)} - p_{ab}^*$

- Expanding and Taylor-approximating:

$$q_{ab}^{(l+1)} = q_{ab}^* + \epsilon_{ab}^{(l+1)} = \sigma_w^2 \mathcal{T}\left(q_{ab}^{(l)}\right) + \sigma_b^2 = \sigma_w^2 \mathcal{T}\left(q_{ab}^* + \epsilon_{ab}^{(l)}\right) + \sigma_b^2$$

$$\overset{\textcolor{red}{Taylor}}{=} \sigma_w^2 \mathcal{T}(q_{ab}^*) + \sigma_b^2 + \sigma_w^2 \dot{\mathcal{T}}(q_{ab}^*)\, \epsilon_{ab}^{(l)} + \mathcal{O}\left(\left(\epsilon_{ab}^{(l)}\right)^2\right)$$

$$= q_{ab}^* + \sigma_w^2 \dot{\mathcal{T}}(q_{ab}^*)\, \epsilon_{ab}^{(l)} + \mathcal{O}\left(\left(\epsilon_{ab}^{(l)}\right)^2\right)$$

- **Hence:** $\epsilon_{ab}^{(l+1)} \approx \sigma_w^2 \dot{\mathcal{T}}(q_{ab}^*)\, \epsilon_{ab}^{(l)} = f'(c^*)^{l+1} \epsilon_{ab}^{(0)}$

# Finite Depth Correction

- Define $\epsilon_{ab}^{(l)} = q_{ab}^{(l)} - q_{ab}^*$ and $\delta_{ab}^{(l)} = p_{ab}^{(l)} - p_{ab}^*$

- Expanding and Taylor-approximating:

$$q_{ab}^{(l+1)} = q_{ab}^* + \epsilon_{ab}^{(l+1)} = \sigma_w^2 \mathcal{T}\left(q_{ab}^{(l)}\right) + \sigma_b^2 = \sigma_w^2 \mathcal{T}\left(q_{ab}^* + \epsilon_{ab}^{(l)}\right) + \sigma_b^2$$

$$\overset{Taylor}{=} \sigma_w^2 \mathcal{T}(q_{ab}^*) + \sigma_b^2 + \sigma_w^2 \dot{\mathcal{T}}(q_{ab}^*) \epsilon_{ab}^{(l)} + \mathcal{O}\left(\left(\epsilon_{ab}^{(l)}\right)^2\right)$$

$$= q_{ab}^* + \sigma_w^2 \dot{\mathcal{T}}(q_{ab}^*) \epsilon_{ab}^{(l)} + \mathcal{O}\left(\left(\epsilon_{ab}^{(l)}\right)^2\right)$$

- **Hence:** $\epsilon_{ab}^{(l+1)} \approx \sigma_w^2 \dot{\mathcal{T}}(q_{ab}^*) \epsilon_{ab}^{(l)} = f'(c^*)^{l+1} \epsilon_{ab}^{(0)}$

- **Similarly:** $\delta_{ab}^{(l+1)} \approx f'(c^*)^{l+1} \left(\delta_{ab}^{(0)} + l\left(1 + \frac{f''(c^*)}{f'(c^*)} p_{ab}^*\right) \epsilon_{ab}^{(0)}\right)$

# Chaotic Regime: Generalization

## Chaotic Regime: Generalization

- Let's analyze the predictive function in the limit:

# Chaotic Regime: Generalization

- Let's analyze the predictive function in the limit:

$$P(\mathbf{\Theta}^{(l)})\mathbf{y} = \left(\mathbf{\Theta}^{(l)}(\mathbf{x}, \mathbf{X})\right)^T \left(\mathbf{\Theta}^{(l)}(\mathbf{X}, \mathbf{X})\right)^{-1} \mathbf{y}$$

## Chaotic Regime: Generalization

- Let's analyze the predictive function in the limit:

$$P(\Theta^{(l)})\boldsymbol{y} = \left(\Theta^{(l)}(\boldsymbol{x}, \boldsymbol{X})\right)^T \left(\Theta^{(l)}(\boldsymbol{X}, \boldsymbol{X})\right)^{-1} \boldsymbol{y}$$
$$\approx (p_{ab}^* \boldsymbol{1})^T \left(p_{ab}^*(\boldsymbol{1}\boldsymbol{1}^T - \mathbb{1}) + p_{aa}^{(l)}\mathbb{1}\right)^{-1} \boldsymbol{y}$$

## Chaotic Regime: Generalization

- Let's analyze the predictive function in the limit:

$$P(\Theta^{(l)})\boldsymbol{y} = \left(\Theta^{(l)}(\boldsymbol{x}, \boldsymbol{X})\right)^T \left(\Theta^{(l)}(\boldsymbol{X}, \boldsymbol{X})\right)^{-1} \boldsymbol{y}$$

$$\approx (p_{ab}^* \mathbf{1})^T \left(p_{ab}^*(\mathbf{1}\mathbf{1}^T - \mathbb{1}) + p_{aa}^{(l)}\mathbb{1}\right)^{-1} \boldsymbol{y}$$

$$= \frac{1}{p_{aa}^{(l)}}(p_{ab}^* \mathbf{1})^T \left(\frac{p_{ab}^*}{p_{aa}^{(l)}}(\mathbf{1}\mathbf{1}^T - \mathbb{1}) + \mathbb{1}\right)^{-1} \boldsymbol{y}$$

# Chaotic Regime: Generalization

- Let's analyze the predictive function in the limit:

$$
\begin{aligned}
P(\boldsymbol{\Theta}^{(l)})\boldsymbol{y} &= \left(\boldsymbol{\Theta}^{(l)}(\boldsymbol{x}, \boldsymbol{X})\right)^{T} \left(\boldsymbol{\Theta}^{(l)}\left(\boldsymbol{X}, \boldsymbol{X}\right)\right)^{-1} \boldsymbol{y} \\
&\approx (p_{ab}^{*}\boldsymbol{1})^{T} \left(p_{ab}^{*}(\boldsymbol{1}\boldsymbol{1}^{T} - \mathbb{1}) + p_{aa}^{(l)}\mathbb{1}\right)^{-1} \boldsymbol{y} \\
&= \frac{1}{p_{aa}^{(l)}}(p_{ab}^{*}\boldsymbol{1})^{T} \left(\frac{p_{ab}^{*}}{p_{aa}^{(l)}}(\boldsymbol{1}\boldsymbol{1}^{T} - \mathbb{1}) + \mathbb{1}\right)^{-1} \boldsymbol{y} \\
&\xrightarrow{l \to \infty} \boldsymbol{0}
\end{aligned}
$$

## Chaotic Regime: Generalization

- Let's analyze the predictive function in the limit:

$$P(\Theta^{(l)})\boldsymbol{y} = \left(\Theta^{(l)}(\boldsymbol{x}, \boldsymbol{X})\right)^T \left(\Theta^{(l)}(\boldsymbol{X}, \boldsymbol{X})\right)^{-1} \boldsymbol{y}$$

$$\approx (p_{ab}^* \mathbf{1})^T \left(p_{ab}^*(\mathbf{11}^T - \mathbb{1}) + p_{aa}^{(l)}\mathbb{1}\right)^{-1} \boldsymbol{y}$$

$$= \frac{1}{p_{aa}^{(l)}}(p_{ab}^* \mathbf{1})^T \left(\frac{p_{ab}^*}{p_{aa}^{(l)}}(\mathbf{11}^T - \mathbb{1}) + \mathbb{1}\right)^{-1} \boldsymbol{y}$$

$$\xrightarrow{l \to \infty} \mathbf{0}$$

- We hence get the trivial prediction **independent** of the data!

# Chaotic Regime: Summary

## Chaotic Regime: Summary

- Characterized by $f'(1) = \sigma_w^2 \int \left[ \phi'(\sqrt{q^*}z) \right]^2 \mathcal{D}z > 1$

## Chaotic Regime: Summary

- Characterized by $f'(1) = \sigma_w^2 \int \left[ \phi'(\sqrt{q^*}z) \right]^2 \mathcal{D}z > 1$

- $\lim_{l \to \infty} \lambda_i \left( \boldsymbol{\Theta}^{(l)}(\boldsymbol{X}, \boldsymbol{X}) \right) = \infty \; \forall i$

## Chaotic Regime: Summary

- Characterized by $f'(1) = \sigma_w^2 \int \left[\phi'(\sqrt{q^*}z)\right]^2 \mathcal{D}z > 1$

- $\lim_{l \to \infty} \lambda_i \left(\mathbf{\Theta}^{(l)}(\mathbf{X}, \mathbf{X})\right) = \infty \ \forall i$

- $\lim_{l \to \infty} \kappa \left(\mathbf{\Theta}^{(l)}(\mathbf{X}, \mathbf{X})\right) = 1 \implies$ Remains **trainable**

## Chaotic Regime: Summary

- Characterized by $f'(1) = \sigma_w^2 \int \left[\phi'(\sqrt{q^*}z)\right]^2 \mathcal{D}z > 1$

- $\lim_{l\to\infty} \lambda_i \left(\boldsymbol{\Theta}^{(l)}(\boldsymbol{X}, \boldsymbol{X})\right) = \infty \;\; \forall i$

- $\lim_{l\to\infty} \kappa \left(\boldsymbol{\Theta}^{(l)}(\boldsymbol{X}, \boldsymbol{X})\right) = 1 \implies$ Remains **trainable**

- $\lim_{l\to\infty} f_\infty(\boldsymbol{x}) = \lim_{l\to\infty} P\left(\boldsymbol{\Theta}^{(l)}(\boldsymbol{x}, \boldsymbol{X})\right)\boldsymbol{y} = \boldsymbol{0} \;\; \forall \boldsymbol{x} \notin \boldsymbol{X}$

# Chaotic Regime: Summary

- Characterized by $f'(1) = \sigma_w^2 \int \left[\phi'(\sqrt{q^*}z)\right]^2 \mathcal{D}z > 1$

- $\lim_{l\to\infty} \lambda_i\left(\boldsymbol{\Theta}^{(l)}(\boldsymbol{X}, \boldsymbol{X})\right) = \infty \ \forall i$

- $\lim_{l\to\infty} \kappa\left(\boldsymbol{\Theta}^{(l)}(\boldsymbol{X}, \boldsymbol{X})\right) = 1 \implies$ Remains **trainable**

- $\lim_{l\to\infty} f_\infty(\boldsymbol{x}) = \lim_{l\to\infty} P\left(\boldsymbol{\Theta}^{(l)}(\boldsymbol{x}, \boldsymbol{X})\right)\boldsymbol{y} = \boldsymbol{0} \quad \forall \boldsymbol{x} \notin \boldsymbol{X}$

    $\implies$ Network **fails** to generalize

**Ordered Regime: Generalization (I)**

## Ordered Regime: Generalization (I)

- More **subtle**: Can't use continuity of map

$$f : GL_n(\mathbb{R}) \to GL_n(\mathbb{R}), \ \boldsymbol{A} \mapsto f(\boldsymbol{A}) = \boldsymbol{A}^{-1}$$

because $\det\left(\boldsymbol{\Theta}^{(\infty)}\right) = 0$ (can't just plugin the limit)

## Ordered Regime: Generalization (I)

- More **subtle**: Can't use continuity of map

$$f : GL_n(\mathbb{R}) \to GL_n(\mathbb{R}), \ \ \boldsymbol{A} \mapsto f(\boldsymbol{A}) = \boldsymbol{A}^{-1}$$

  because $\det\left(\boldsymbol{\Theta}^{(\infty)}\right) = 0$ (can't just plugin the limit)

- Have to work with the finite depth approximation to argue like

$$\lim_{\lambda \to \infty} (\boldsymbol{A} + \lambda \mathbb{1})^{-1} = \boldsymbol{0}$$

  to obtain the correct limit

# Ordered Regime: Generalization (I)

- More **subtle**: Can't use continuity of map

$$f : GL_n(\mathbb{R}) \to GL_n(\mathbb{R}), \quad \boldsymbol{A} \mapsto f(\boldsymbol{A}) = \boldsymbol{A}^{-1}$$

because $\det\left(\boldsymbol{\Theta}^{(\infty)}\right) = 0$ (can't just plugin the limit)

- Have to work with the finite depth approximation to argue like

$$\lim_{\lambda \to \infty} (\boldsymbol{A} + \lambda \mathbb{1})^{-1} = \boldsymbol{0}$$

to obtain the correct limit

- Write $\boldsymbol{\Theta}^{(l)}(\boldsymbol{X}, \boldsymbol{X}) = p^* \boldsymbol{11}^T + l \left(f'(1)\right)^l \boldsymbol{A}^{(l)}(\boldsymbol{X}, \boldsymbol{X})$

*Handwritten annotation:*
- Width correction matrix that depends on the data
- $A^{(l)}(X,X) \to A^{(\infty)}(X,X)$

# Ordered Regime: Generalization (II)

## Ordered Regime: Generalization (II)

Let $\lambda = jf'(1)^j$, $\boldsymbol{A}^{(j)}(\boldsymbol{X}, \boldsymbol{X}) = \boldsymbol{A}_{\boldsymbol{X}}^{(j)}$. Let's express the finite depth inverse:

# Ordered Regime: Generalization (II)

Let $\lambda = jf'(1)^j$, $\boldsymbol{A}^{(j)}(\boldsymbol{X}, \boldsymbol{X}) = \boldsymbol{A}_{\boldsymbol{X}}^{(j)}$. Let's express the finite depth inverse:

$$\left(\boldsymbol{\Theta}^{(j)}(\boldsymbol{X}, \boldsymbol{X})\right)^{-1} = \left(p^* \mathbf{1}\mathbf{1}^T + \lambda \boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1}$$

## Ordered Regime: Generalization (II)

Let $\lambda = jf'(1)^j$, $\boldsymbol{A}^{(j)}(\boldsymbol{X}, \boldsymbol{X}) = \boldsymbol{A}_{\boldsymbol{X}}^{(j)}$. Let's express the finite depth inverse:

$$\left(\boldsymbol{\Theta}^{(j)}(\boldsymbol{X}, \boldsymbol{X})\right)^{-1} = \left(p^* \boldsymbol{1}\boldsymbol{1}^T + \lambda \boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1}$$

$(A + uV^T)^{-1} = A^{-1} + \dfrac{A^{-1} u V^T A^{-1}}{1 + V^T A^{-1} u}$

$$= \lambda^{-1}\left(\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1} - \frac{\lambda^{-2}\left(\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1}\boldsymbol{1}\boldsymbol{1}^T\left(\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1}}{\frac{1}{p^*} + \lambda^{-1}\boldsymbol{1}^T\left(\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1}\boldsymbol{1}}$$

# Ordered Regime: Generalization (II)

Let $\lambda = jf'(1)^j$, $\boldsymbol{A}^{(j)}(\boldsymbol{X}, \boldsymbol{X}) = \boldsymbol{A}_{\boldsymbol{X}}^{(j)}$. Let's express the finite depth inverse:

$$\left(\boldsymbol{\Theta}^{(j)}(\boldsymbol{X}, \boldsymbol{X})\right)^{-1} = \left(p^* \mathbf{1}\mathbf{1}^T + \lambda \boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1}$$

$$(A + uv^T)^{-1} = A^{-1} + \frac{A^{-1} u v^T A^{-1}}{1 + v^T A^{-1} u}$$

$$= \lambda^{-1} \left(\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1} - \frac{\lambda^{-2} \left(\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1} \mathbf{1}\mathbf{1}^T \left(\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1}}{\frac{1}{p^*} + \lambda^{-1} \mathbf{1}^T \left(\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1} \mathbf{1}}$$

$$= \lambda^{-1} \left(\left(\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1} - \frac{\left(\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1} \mathbf{1}\mathbf{1}^T \left(\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1}}{\lambda \left(\frac{1}{p^*} + \lambda^{-1} \mathbf{1}^T \left(\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1} \mathbf{1}\right)}\right)$$

# Ordered Regime: Generalization (II)

Let $\lambda = jf'(1)^j$, $\boldsymbol{A}^{(j)}(\boldsymbol{X}, \boldsymbol{X}) = \boldsymbol{A}_{\boldsymbol{X}}^{(j)}$. Let's express the finite depth inverse:

$$\left(\boldsymbol{\Theta}^{(j)}(\boldsymbol{X}, \boldsymbol{X})\right)^{-1} = \left(p^*\mathbf{1}\mathbf{1}^T + \lambda\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1}$$

$(A + uv^T)^{-1} = A^{-1} + \dfrac{A^{-1}uv^TA^{-1}}{1 + v^TA^{-1}u}$

$$= \lambda^{-1}\left(\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1} - \frac{\lambda^{-2}\left(\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1}\mathbf{1}\mathbf{1}^T\left(\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1}}{\frac{1}{p^*} + \lambda^{-1}\mathbf{1}^T\left(\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1}\mathbf{1}}$$

$$= \lambda^{-1}\left(\left(\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1} - \frac{\left(\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1}\mathbf{1}\mathbf{1}^T\left(\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1}}{\lambda\left(\frac{1}{p^*} + \lambda^{-1}\mathbf{1}^T\left(\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1}\mathbf{1}\right)}\right)$$

$$= \lambda^{-1}\left(\left(\boldsymbol{A}_{\boldsymbol{X}}^{(j)}\right)^{-1} - \hat{p}\boldsymbol{a}\boldsymbol{a}^T\right)$$

$a = \left(A_{\boldsymbol{X}}^{(j)}\right)^{-1}\underline{1}$

# Ordered Regime: Generalization (III)

# Ordered Regime: Generalization (III)

- Take a new input $z \implies \Theta^{(j)}(z, X) = p^* \mathbf{1}^T + \lambda A_z^{(j)}$

## Ordered Regime: Generalization (III)

- Take a new input $z \implies \Theta^{(j)}(z, X) = p^* \mathbf{1}^T + \lambda A_z^{(j)}$

- Let's express the predictive function:

## Ordered Regime: Generalization (III)

- Take a new input $z \implies \Theta^{(j)}(z, X) = p^* 1^T + \lambda A_z^{(j)}$

- Let's express the predictive function:

$$P(\Theta^{(j)}) = \left( p^* 1^T + \lambda A_z^{(j)} \right) \lambda^{-1} \left( \left( A_X^{(j)} \right)^{-1} - \hat{\rho} a a^T \right)$$

## Ordered Regime: Generalization (III)

- Take a new input $\mathbf{z} \implies \mathbf{\Theta}^{(j)}(\mathbf{z}, \mathbf{X}) = p^* \mathbf{1}^T + \lambda \mathbf{A}_{\mathbf{z}}^{(j)}$

- Let's express the predictive function:

$$P(\mathbf{\Theta}^{(j)}) = \left( p^* \mathbf{1}^T + \lambda \mathbf{A}_{\mathbf{z}}^{(j)} \right) \lambda^{-1} \left( \left( \mathbf{A}_{\mathbf{X}}^{(j)} \right)^{-1} - \hat{p} \mathbf{a} \mathbf{a}^T \right)$$

$$= \mathbf{A}_{\mathbf{z}}^{(j)} \left( \mathbf{A}_{\mathbf{X}}^{(j)} \right)^{-1} - \hat{p} \mathbf{A}_{\mathbf{z}}^{(j)} \mathbf{a} \mathbf{a}^T + \lambda^{-1} p^* \left( \mathbf{a}^T - \hat{p} \mathbf{1}^T \mathbf{a} \mathbf{a}^T \right)$$

## Ordered Regime: Generalization (III)

- Take a new input $z \implies \Theta^{(j)}(z, X) = p^* \mathbf{1}^T + \lambda A_z^{(j)}$

- Let's express the predictive function:

$$
\begin{aligned}
P(\Theta^{(j)}) &= \left( p^* \mathbf{1}^T + \lambda A_z^{(j)} \right) \lambda^{-1} \left( \left( A_X^{(j)} \right)^{-1} - \hat{p} a a^T \right) \\
&= A_z^{(j)} \left( A_X^{(j)} \right)^{-1} - \hat{p} A_z^{(j)} a a^T + \lambda^{-1} p^* \left( a^T - \hat{p} \mathbf{1}^T a a^T \right) \\
&= A_z^{(j)} \left( A_X^{(j)} \right)^{-1} - \hat{p} A_z^{(j)} a a^T + \hat{p} a^T
\end{aligned}
$$

## Ordered Regime: Generalization (III)

- Take a new input $z \implies \Theta^{(j)}(z, X) = p^* 1^T + \lambda A_z^{(j)}$

- Let's express the predictive function:

$$\begin{aligned}
P(\Theta^{(j)}) &= \left( p^* 1^T + \lambda A_z^{(j)} \right) \lambda^{-1} \left( \left( A_X^{(j)} \right)^{-1} - \hat{p} a a^T \right) \\
&= A_z^{(j)} \left( A_X^{(j)} \right)^{-1} - \hat{p} A_z^{(j)} a a^T + \lambda^{-1} p^* \left( a^T - \hat{p} 1^T a a^T \right) \\
&= A_z^{(j)} \left( A_X^{(j)} \right)^{-1} - \hat{p} A_z^{(j)} a a^T + \hat{p} a^T \\
&\xrightarrow{j \to \infty} A_z^{(\infty)} \left( A_X^{(\infty)} \right)^{-1} - \hat{p} A_z^{(\infty)} a^{(\infty)} \left( a^{(\infty)} \right)^T + \hat{p} \left( a^{(\infty)} \right)^T
\end{aligned}$$

## Ordered Regime: Generalization (III)

- Take a new input $z \implies \Theta^{(j)}(z, X) = p^* 1^T + \lambda A_z^{(j)}$

- Let's express the predictive function:

$$
\begin{aligned}
P(\Theta^{(j)}) &= \left( p^* 1^T + \lambda A_z^{(j)} \right) \lambda^{-1} \left( \left( A_X^{(j)} \right)^{-1} - \hat{p} a a^T \right) \\
&= A_z^{(j)} \left( A_X^{(j)} \right)^{-1} - \hat{p} A_z^{(j)} a a^T + \lambda^{-1} p^* \left( a^T - \hat{p} 1^T a a^T \right) \\
&= A_z^{(j)} \left( A_X^{(j)} \right)^{-1} - \hat{p} A_z^{(j)} a a^T + \hat{p} a^T \\
&\xrightarrow{j \to \infty} A_z^{(\infty)} \left( A_X^{(\infty)} \right)^{-1} - \hat{p} A_z^{(\infty)} a^{(\infty)} \left( a^{(\infty)} \right)^T + \hat{p} \left( a^{(\infty)} \right)^T
\end{aligned}
$$

- **Thus:** $f_\infty(z) = \left( A_z^{(\infty)} \left( A_X^{(\infty)} \right)^{-1} + \hat{A} \right) y$

## Ordered Regime: Generalization (III)

- Take a new input $\boldsymbol{z} \implies \Theta^{(j)}(\boldsymbol{z}, \boldsymbol{X}) = p^* \mathbf{1}^T + \lambda \boldsymbol{A}_{\boldsymbol{z}}^{(j)}$

- Let's express the predictive function:

$$
\begin{aligned}
P(\Theta^{(j)}) &= \left( p^* \mathbf{1}^T + \lambda \boldsymbol{A}_{\boldsymbol{z}}^{(j)} \right) \lambda^{-1} \left( \left( \boldsymbol{A}_{\boldsymbol{X}}^{(j)} \right)^{-1} - \hat{p} \boldsymbol{a} \boldsymbol{a}^T \right) \\
&= \boldsymbol{A}_{\boldsymbol{z}}^{(j)} \left( \boldsymbol{A}_{\boldsymbol{X}}^{(j)} \right)^{-1} - \hat{p} \boldsymbol{A}_{\boldsymbol{z}}^{(j)} \boldsymbol{a} \boldsymbol{a}^T + \lambda^{-1} p^* \left( \boldsymbol{a}^T - \hat{p} \mathbf{1}^T \boldsymbol{a} \boldsymbol{a}^T \right) \\
&= \boldsymbol{A}_{\boldsymbol{z}}^{(j)} \left( \boldsymbol{A}_{\boldsymbol{X}}^{(j)} \right)^{-1} - \hat{p} \boldsymbol{A}_{\boldsymbol{z}}^{(j)} \boldsymbol{a} \boldsymbol{a}^T + \hat{p} \boldsymbol{a}^T \\
&\xrightarrow{j \to \infty} \boldsymbol{A}_{\boldsymbol{z}}^{(\infty)} \left( \boldsymbol{A}_{\boldsymbol{X}}^{(\infty)} \right)^{-1} - \hat{p} \boldsymbol{A}_{\boldsymbol{z}}^{(\infty)} \boldsymbol{a}^{(\infty)} \left( \boldsymbol{a}^{(\infty)} \right)^T + \hat{p} \left( \boldsymbol{a}^{(\infty)} \right)^T
\end{aligned}
$$

- **Thus:** $f_\infty(\boldsymbol{z}) = \left( \boldsymbol{A}_{\boldsymbol{z}}^{(\infty)} \left( \boldsymbol{A}_{\boldsymbol{X}}^{(\infty)} \right)^{-1} + \hat{\boldsymbol{A}} \right) \boldsymbol{y}$

  **Non-trivial** generalization possible!

# Ordered Regime: Summary

# Ordered Regime: Summary

- Characterized by $f'(1) = \sigma_w^2 \int \left[\phi'(\sqrt{q^*}z)\right]^2 \mathcal{D}z < 1$

## Ordered Regime: Summary

- Characterized by $f'(1) = \sigma_w^2 \int \left[\phi'(\sqrt{q^*}z)\right]^2 \mathcal{D}z < 1$

- $\lim_{l\to\infty} \lambda_1\left(\Theta^{(l)}(\boldsymbol{X}, \boldsymbol{X})\right) = np^*$

# Ordered Regime: Summary

- Characterized by $f'(1) = \sigma_w^2 \int \left[\phi'(\sqrt{q^*}z)\right]^2 \mathcal{D}z < 1$

- $\lim_{l\to\infty} \lambda_1 \left(\Theta^{(l)}(\boldsymbol{X},\boldsymbol{X})\right) = np^*$

- $\lim_{l\to\infty} \lambda_{\text{bulk}} \left(\Theta^{(l)}(\boldsymbol{X},\boldsymbol{X})\right) = 0$
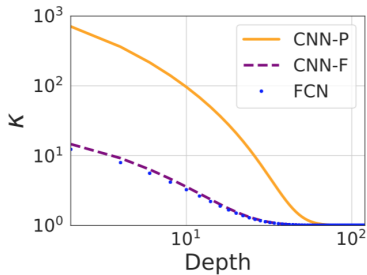
## Ordered Regime: Summary

- Characterized by $f'(1) = \sigma_w^2 \int \left[ \phi'(\sqrt{q^*}z) \right]^2 \mathcal{D}z < 1$

- $\lim_{l \to \infty} \lambda_1 \left( \Theta^{(l)}(\boldsymbol{X}, \boldsymbol{X}) \right) = np^*$

- $\lim_{l \to \infty} \lambda_{\text{bulk}} \left( \Theta^{(l)}(\boldsymbol{X}, \boldsymbol{X}) \right) = 0$

- $\lim_{l \to \infty} \kappa \left( \Theta^{(l)}(\boldsymbol{X}, \boldsymbol{X}) \right) = \infty \implies$ Becomes **untrainable**
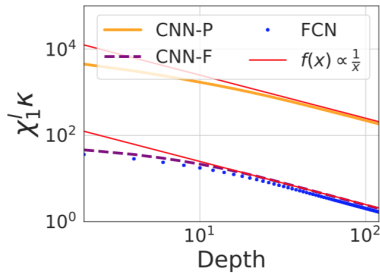
## Ordered Regime: Summary

- Characterized by $f'(1) = \sigma_w^2 \int \left[\phi'(\sqrt{q^*}z)\right]^2 \mathcal{D}z < 1$

- $\lim_{l \to \infty} \lambda_1 \left(\Theta^{(l)}(\boldsymbol{X}, \boldsymbol{X})\right) = np^*$

- $\lim_{l \to \infty} \lambda_{\text{bulk}} \left(\Theta^{(l)}(\boldsymbol{X}, \boldsymbol{X})\right) = 0$

- $\lim_{l \to \infty} \kappa \left(\Theta^{(l)}(\boldsymbol{X}, \boldsymbol{X})\right) = \infty \implies$ Becomes **untrainable**

- $\lim_{l \to \infty} f_\infty^{(l)}(\boldsymbol{x}) = \left(\boldsymbol{A}_z^{(\infty)} \left(\boldsymbol{A}_{\boldsymbol{X}}^{(\infty)}\right)^{-1} + \hat{\boldsymbol{A}}\right) \boldsymbol{y}$

# Ordered Regime: Summary

- Characterized by $f'(1) = \sigma_w^2 \int \left[\phi'(\sqrt{q^*}z)\right]^2 \mathcal{D}z < 1$

- $\lim_{l \to \infty} \lambda_1\left(\Theta^{(l)}(\boldsymbol{X}, \boldsymbol{X})\right) = np^*$

- $\lim_{l \to \infty} \lambda_{\text{bulk}}\left(\Theta^{(l)}(\boldsymbol{X}, \boldsymbol{X})\right) = 0$

- $\lim_{l \to \infty} \kappa\left(\Theta^{(l)}(\boldsymbol{X}, \boldsymbol{X})\right) = \infty \implies$ Becomes **untrainable**

- $\lim_{l \to \infty} f_\infty^{(l)}(\boldsymbol{x}) = \left(\boldsymbol{A}_z^{(\infty)}\left(\boldsymbol{A}_{\boldsymbol{X}}^{(\infty)}\right)^{-1} + \hat{\boldsymbol{A}}\right)\boldsymbol{y}$

  $\implies$ Network might be **able** to generalize

# Experiments (I)



Chaotic: $K \to 1$

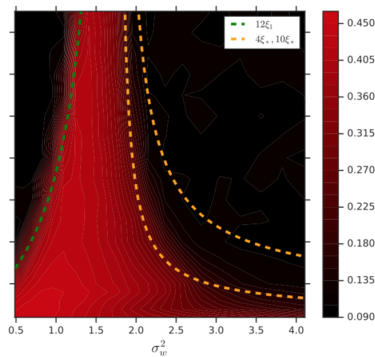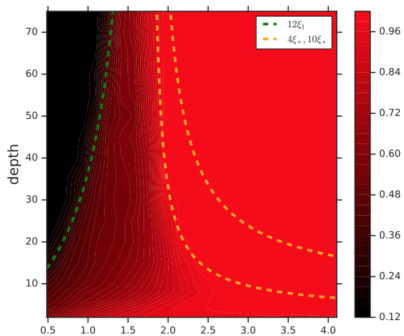ordered: $K \approx (\chi_1^l \cdot l)^{-1} \to \infty$

Condition number in the chaotic and ordered regime

# Experiments (II)



Train accuracy

Test accuracy

Test and Training Accuracy for Finite Depth and Width Networks

## Discussion

- Important to also study **infinite depth**, after all seems to be more important in practice than big widths

- Interesting **tradeoff** between generalization and trainability

- Characterization of generalization a bit **weak** but seems to be true empirically

- Approximation of the diagonal $\Sigma^{(l)}(\boldsymbol{x}, \boldsymbol{x}) = q^*$ **very strong** assumption, unclear how this affects the obtained limits

- Very **messy** paper, lots of referencing to prior work and some statements are a bit unclear