

Lecture 3:

The Principle of Sparsity

[§4, §8, §11.2, and § ∞ .3 of

“The Principles of Deep Learning Theory (PDLT),”
arXiv:2106.10165]

Problems 1, 2, & 3

Fully-trained network output, Taylor-expanded around initialization:

$$z^{\star} = z + (\theta^{\star} - \theta) \frac{dz}{d\theta} + \frac{1}{2} (\theta^{\star} - \theta)^2 \frac{d^2 z}{d\theta^2} + \dots$$

Problems 1, 2, & 3

Fully-trained network output, Taylor-expanded around initialization:

$$z^{\star} = z + (\theta^{\star} - \theta) \frac{dz}{d\theta} + \frac{1}{2} (\theta^{\star} - \theta)^2 \frac{d^2 z}{d\theta^2} + \dots$$

- Problem 1: too many terms in general

Problems 1, 2, & 3

Fully-trained network output, Taylor-expanded around initialization:

$$z^{\star} = z + (\theta^{\star} - \theta) \frac{dz}{d\theta} + \frac{1}{2} (\theta^{\star} - \theta)^2 \frac{d^2 z}{d\theta^2} + \dots$$

- Problem 1: too many terms in general
- Problem 2: complicated mapping

$$p(\theta) \rightarrow p \left(\theta, z, \frac{dz}{d\theta}, \frac{d^2 z}{d\theta^2}, \dots \right)$$

initial distributions
over
model parameters

statistics at *initialization*

Problems 1, 2, & 3

Fully-trained network output, Taylor-expanded around initialization:

$$z^{\star} = z + (\theta^{\star} - \theta) \frac{dz}{d\theta} + \frac{1}{2} (\theta^{\star} - \theta)^2 \frac{d^2 z}{d\theta^2} + \dots$$

- Problem 1: too many terms in general
- Problem 2: complicated mapping
- Problem 3: complicated dynamics $\theta^{\star} = [\theta^{\star}] \left(\theta, z, \frac{dz}{d\theta}, \frac{d^2 z}{d\theta^2}, \dots; \text{algorithm}; \text{data} \right)$

$$p(\theta) \rightarrow p \left(\theta, z, \frac{dz}{d\theta}, \frac{d^2 z}{d\theta^2}, \dots \right)$$

initial distributions
over
model parameters

statistics at *initialization*

Dan has covered Dynamics

Fully-trained network output, Taylor-expanded around initialization:

$$z^{\star} = z + (\theta^{\star} - \theta) \frac{dz}{d\theta} + \frac{1}{2} (\theta^{\star} - \theta)^2 \frac{d^2 z}{d\theta^2} + \dots$$

- ~~Problem 1: too many terms in general~~

- Problem 2: complicated mapping

- **Problem 3: complicated dynamics**

$$z^{\star} = [z^{\star}] \left(\theta, z, \frac{dz}{d\theta}, \frac{d^2 z}{d\theta^2}, \dots; \text{algorithm}; \text{data} \right)$$

$$p(\theta) \rightarrow p \left(\theta, z, \frac{dz}{d\theta}, \frac{d^2 z}{d\theta^2}, \dots \right)$$

initial distributions
over
model parameters

statistics at *initialization*

Dan has covered Dynamics

Fully-trained network output, Taylor-expanded around initialization:

$$z^{\star} = z + (\theta^{\star} - \theta) \frac{dz}{d\theta} + \frac{1}{2} (\theta^{\star} - \theta)^2 \frac{d^2 z}{d\theta^2} + \dots$$

- ~~Problem 1: too many terms in general~~

- Problem 2: complicated mapping

- **Problem 3: complicated dynamics**

$$z^{\star} = [z^{\star}] \left(\theta, z, \frac{dz}{d\theta}, \frac{d^2 z}{d\theta^2}, \dots; \text{algorithm; data} \right)$$

$$p(\theta) \rightarrow p \left(\theta, z, \frac{dz}{d\theta}, \frac{d^2 z}{d\theta^2}, \dots \right) \rightarrow p(z^{\star})$$

initial distributions
over
model parameters

statistics at *initialization*

statistics *after training*

Sho will cover Statistics

Fully-trained network output, Taylor-expanded around initialization:

$$z^{\star} = z + (\theta^{\star} - \theta) \frac{dz}{d\theta} + \frac{1}{2} (\theta^{\star} - \theta)^2 \frac{d^2 z}{d\theta^2} + \dots$$

- Problem 1: too many terms in general
- Problem 2: complicated mapping
- Problem 3: complicated dynamics

$$p(\theta) \rightarrow p\left(\theta, z, \frac{dz}{d\theta}, \frac{d^2 z}{d\theta^2}, \dots\right) \rightarrow p(z^{\star})$$

initial distributions
over
model parameters

statistics at *initialization*

statistics *after training*

Sho will cover Statistics

$$p(\theta) \rightarrow p\left(\theta, z, \frac{dz}{d\theta}, \frac{d^2z}{d\theta^2}, \dots\right)$$

initial distributions
over
model parameters

statistics at *initialization*
for WIDE & DEEP neural networks

Sho will cover Statistics

Lecture 3: The Principle of Sparsity, deriving recursions

Lecture 4: The Principle of Criticality, solving recursions

$$p(\theta) \rightarrow p\left(\theta, z, \frac{dz}{d\theta}, \frac{d^2z}{d\theta^2}, \dots\right)$$

initial distributions
over
model parameters

statistics at *initialization*
for WIDE & DEEP neural networks

Outline

1. Neural Networks 101
2. One-Layer Neural Networks
3. Two-Layer Neural Networks
4. Deep Neural Networks

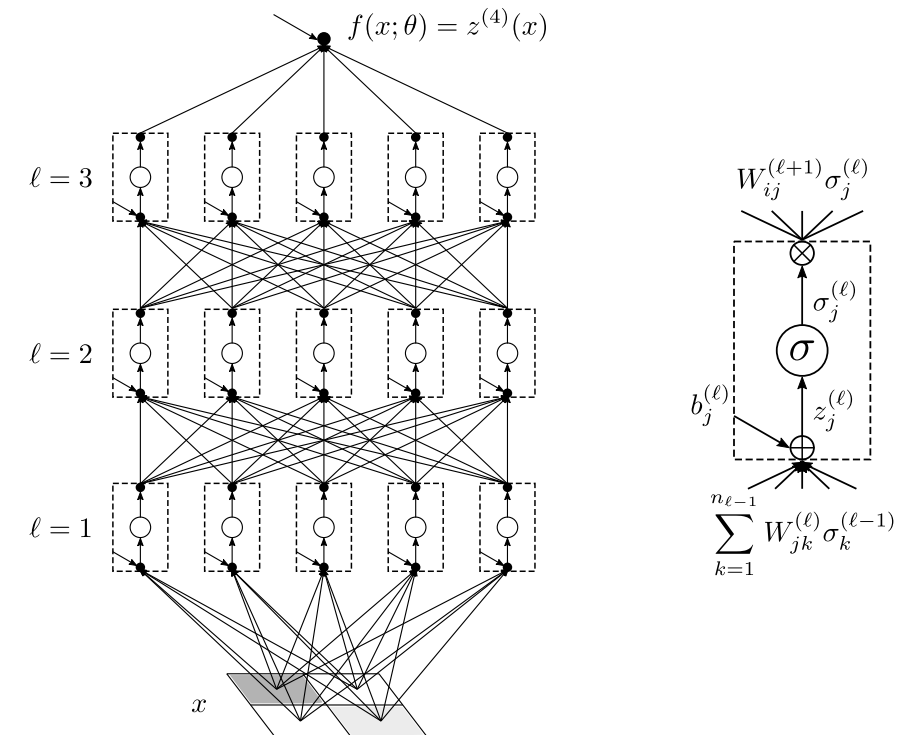
1. Neural Networks 101

Neural Networks

$$\hat{z}_i^{(1)}(x) \equiv b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j \quad \text{for } i = 1, \dots, n_1,$$

$$\hat{z}_i^{(\ell+1)}(x) \equiv b_i^{(\ell+1)} + \sum_{j=1}^{n_\ell} W_{ij}^{(\ell+1)} \sigma\left(\hat{z}_j^{(\ell)}(x)\right) \quad \text{for } i = 1, \dots, n_{\ell+1}; \ell = 1, \dots, L-1$$

$$\hat{z}_{i;\delta} = \hat{z}_i^{(L)}(x_\delta)$$



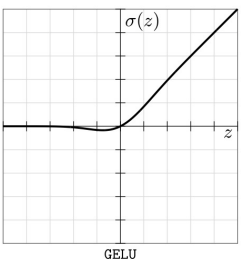
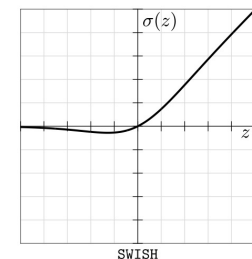
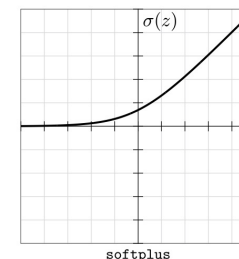
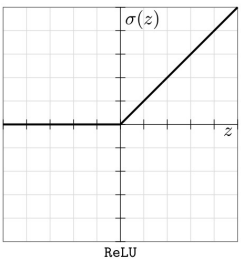
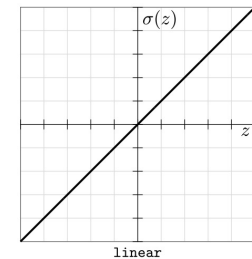
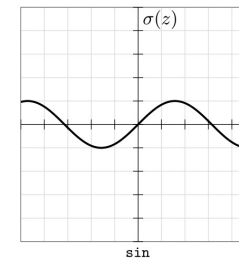
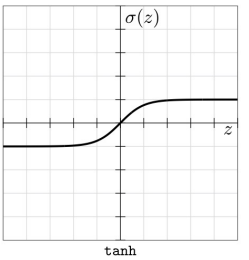
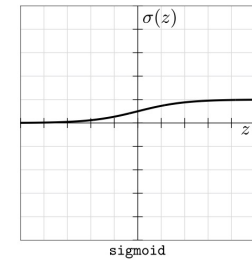
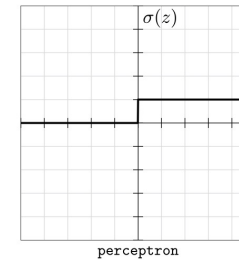
Neural Networks

$$\hat{z}_i^{(1)}(x) \equiv b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j \quad \text{for } i = 1, \dots, n_1,$$

$$\hat{z}_i^{(\ell+1)}(x) \equiv b_i^{(\ell+1)} + \sum_{j=1}^{n_\ell} W_{ij}^{(\ell+1)} \sigma\left(\hat{z}_j^{(\ell)}(x)\right) \quad \text{for } i = 1, \dots, n_{\ell+1}; \ell = 1, \dots, L-1$$

$$\hat{z}_{i;\delta} = \hat{z}_i^{(L)}(x_\delta)$$

activation function $\sigma(z)$



Neural Networks

preactivations

$$\hat{z}_i^{(1)}(x) \equiv b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j \quad \text{for } i = 1, \dots, n_1,$$
$$\hat{z}_i^{(\ell+1)}(x) \equiv b_i^{(\ell+1)} + \sum_{j=1}^{n_\ell} W_{ij}^{(\ell+1)} \sigma\left(\hat{z}_j^{(\ell)}(x)\right) \quad \text{for } i = 1, \dots, n_{\ell+1}; \ell = 1, \dots, L-1$$

network output

$$\hat{z}_{i;\delta} = \hat{z}_i^{(L)}(x_\delta)$$

Neural Networks

$$\hat{z}_i^{(1)}(x) \equiv b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j \quad \text{for } i = 1, \dots, n_1,$$

$$\hat{z}_i^{(\ell+1)}(x) \equiv b_i^{(\ell+1)} + \sum_{j=1}^{n_\ell} W_{ij}^{(\ell+1)} \sigma\left(\hat{z}_j^{(\ell)}(x)\right) \quad \text{for } i = 1, \dots, n_{\ell+1}; \ell = 1, \dots, L-1$$

$$\hat{z}_{i;\delta} = \hat{z}_i^{(L)}(x_\delta)$$

hat @ initialization

Neural Networks

$$\widehat{z}_i^{(1)}(x) \equiv b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j \quad \text{for } i = 1, \dots, n_1,$$

$$\widehat{z}_i^{(\ell+1)}(x) \equiv b_i^{(\ell+1)} + \sum_{j=1}^{n_\ell} W_{ij}^{(\ell+1)} \sigma\left(\widehat{z}_j^{(\ell)}(x)\right) \quad \text{for } i = 1, \dots, n_{\ell+1}; \ell = 1, \dots, L-1$$

$$\widehat{z}_{i;\delta} = \widehat{z}_i^{(L)}(x_\delta)$$

Biases and weights (model parameters) are independently (& symmetrically) distributed with variances

$$\mathbb{E} \left[b_{i_1}^{(\ell)} b_{i_2}^{(\ell)} \right] = \delta_{i_1 i_2} C_b^{(\ell)}, \quad \mathbb{E} \left[W_{i_1 j_1}^{(\ell)} W_{i_2 j_2}^{(\ell)} \right] = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{C_W^{(\ell)}}{n_{\ell-1}}$$

Neural Networks

$$\hat{z}_i^{(1)}(x) \equiv b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j \quad \text{for } i = 1, \dots, n_1,$$

$$\hat{z}_i^{(\ell+1)}(x) \equiv b_i^{(\ell+1)} + \sum_{j=1}^{n_\ell} W_{ij}^{(\ell+1)} \sigma\left(\hat{z}_j^{(\ell)}(x)\right) \quad \text{for } i = 1, \dots, n_{\ell+1}; \ell = 1, \dots, L-1$$

$$\hat{z}_{i;\delta} = \hat{z}_i^{(L)}(x_\delta)$$

Biases and weights (model parameters) are independently (& symmetrically) distributed with variances

$$\mathbb{E} \left[b_{i_1}^{(\ell)} b_{i_2}^{(\ell)} \right] = \delta_{i_1 i_2} C_b^{(\ell)}, \quad \mathbb{E} \left[W_{i_1 j_1}^{(\ell)} W_{i_2 j_2}^{(\ell)} \right] = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{C_W^{(\ell)}}{n_{\ell-1}}$$

initialization hyperparameters

Neural Networks

$$\hat{z}_i^{(1)}(x) \equiv b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j \quad \text{for } i = 1, \dots, n_1,$$

$$\hat{z}_i^{(\ell+1)}(x) \equiv b_i^{(\ell+1)} + \sum_{j=1}^{n_\ell} W_{ij}^{(\ell+1)} \sigma\left(\hat{z}_j^{(\ell)}(x)\right) \quad \text{for } i = 1, \dots, n_{\ell+1}; \ell = 1, \dots, L-1$$

$$\hat{z}_{i;\delta} = \hat{z}_i^{(L)}(x_\delta)$$

Biases and weights (model parameters) are independently (& symmetrically) distributed with variances

$$\mathbb{E} \left[b_{i_1}^{(\ell)} b_{i_2}^{(\ell)} \right] = \delta_{i_1 i_2} C_b^{(\ell)}, \quad \mathbb{E} \left[W_{i_1 j_1}^{(\ell)} W_{i_2 j_2}^{(\ell)} \right] = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{C_W^{(\ell)}}{n_{\ell-1}}$$

good wide limit

One Aside on Gradient Descent

$$\theta_{\mu}(t+1) = \theta_{\mu}(t) - \eta \sum_{\nu} \lambda_{\mu\nu} \left(\sum \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \frac{dz_{j;\tilde{\alpha}}}{d\theta_{\nu}} \right)$$

[Cf. Andrea's "S" matrix]

One Aside on Gradient Descent

$$\theta_{\mu}(t+1) = \theta_{\mu}(t) - \eta \sum_{\nu} \lambda_{\mu\nu} \left(\sum \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \frac{dz_{j;\tilde{\alpha}}}{d\theta_{\nu}} \right)$$

Taylor expansion:

$$\begin{aligned} z_{i;\delta}(t+1) = & z_{i;\delta}(t) \\ & - \eta \sum_{j,\tilde{\alpha}} \left(\sum_{\mu,\nu} \lambda_{\mu\nu} \frac{dz_{i;\delta}}{d\theta_{\mu}} \frac{dz_{j;\tilde{\alpha}}}{d\theta_{\nu}} \right) \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \\ & + \dots \end{aligned}$$

One Aside on Gradient Descent

$$\theta_\mu(t+1) = \theta_\mu(t) - \eta \sum_\nu \lambda_{\mu\nu} \left(\sum \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right)$$

Taylor expansion:

$$\begin{aligned} z_{i;\delta}(t+1) = & z_{i;\delta}(t) \quad \text{Neural Tangent Kernel (NTK)} \quad H(t) \quad [\text{Cf. Dan's } k_{\delta\tilde{\alpha}}] \\ & - \eta \sum_{j,\tilde{\alpha}} \left(\sum_{\mu,\nu} \lambda_{\mu\nu} \frac{dz_{i;\delta}}{d\theta_\mu} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right) \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \\ & + \dots \end{aligned}$$

One Aside on Gradient Descent

$$\theta_{\mu}(t+1) = \theta_{\mu}(t) - \eta \sum_{\nu} \lambda_{\mu\nu} \left(\sum \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \frac{dz_{j;\tilde{\alpha}}}{d\theta_{\nu}} \right)$$

Taylor expansion:

$$\begin{aligned} z_{i;\delta}(t+1) = & z_{i;\delta}(t) \quad \text{Neural Tangent Kernel (NTK)} \quad H(t) \quad [\text{Cf. Dan's } k_{ij;\delta\tilde{\alpha}}^{\text{E}}(\theta)] \\ & - \eta \sum_{j,\tilde{\alpha}} \left(\sum_{\mu,\nu} \lambda_{\mu\nu} \frac{dz_{i;\delta}}{d\theta_{\mu}} \frac{dz_{j;\tilde{\alpha}}}{d\theta_{\nu}} \right) \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \quad \text{differential of NTK (dNTK)} \quad dH(t) \quad [\text{Cf. Dan's } \mu_{\delta\tilde{\alpha}_1\tilde{\alpha}_2}] \\ & + \frac{\eta^2}{2} \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2} \left(\sum_{\mu_1,\nu_1,\mu_2,\nu_2} \lambda_{\mu_1\nu_1} \lambda_{\mu_2\nu_2} \frac{d^2 z_{i;\delta}}{d\theta_{\mu_1} d\theta_{\mu_2}} \frac{dz_{j_1;\tilde{\alpha}_1}}{d\theta_{\nu_1}} \frac{dz_{j_2;\tilde{\alpha}_2}}{d\theta_{\nu_2}} \right) \frac{\partial \mathcal{L}}{\partial z_{j_1;\tilde{\alpha}_1}} \frac{\partial \mathcal{L}}{\partial z_{j_2;\tilde{\alpha}_2}} \\ & + \dots \end{aligned}$$

One Aside on Gradient Descent

$$\theta_\mu(t+1) = \theta_\mu(t) - \eta \sum_\nu \lambda_{\mu\nu} \left(\sum \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right)$$

Taylor expansion:

$$\begin{aligned} z_{i;\delta}(t+1) = & z_{i;\delta}(t) \quad \text{Neural Tangent Kernel (NTK)} \quad H(t) \\ & - \eta \sum_{j,\tilde{\alpha}} \left(\sum_{\mu,\nu} \lambda_{\mu\nu} \frac{dz_{i;\delta}}{d\theta_\mu} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right) \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \quad \text{differential of NTK (dNTK)} \quad dH(t) \\ & + \frac{\eta^2}{2} \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2} \left(\sum_{\mu_1,\nu_1,\mu_2,\nu_2} \lambda_{\mu_1\nu_1} \lambda_{\mu_2\nu_2} \frac{d^2 z_{i;\delta}}{d\theta_{\mu_1} d\theta_{\mu_2}} \frac{dz_{j_1;\tilde{\alpha}_1}}{d\theta_{\nu_1}} \frac{dz_{j_2;\tilde{\alpha}_2}}{d\theta_{\nu_2}} \right) \frac{\partial \mathcal{L}}{\partial z_{j_1;\tilde{\alpha}_1}} \frac{\partial \mathcal{L}}{\partial z_{j_2;\tilde{\alpha}_2}} \\ & - \frac{\eta^3}{6} \sum \left(\sum \lambda_{\mu_1\nu_1} \lambda_{\mu_2\nu_2} \lambda_{\mu_3\nu_3} \frac{d^3 z_{i;\delta}}{d\theta_{\mu_1} d\theta_{\mu_2} d\theta_{\mu_3}} \frac{dz_{j_1;\tilde{\alpha}_1}}{d\theta_{\nu_1}} \frac{dz_{j_2;\tilde{\alpha}_2}}{d\theta_{\nu_2}} \frac{dz_{j_3;\tilde{\alpha}_3}}{d\theta_{\nu_3}} \right) \frac{\partial \mathcal{L}}{\partial z_{j_1;\tilde{\alpha}_1}} \frac{\partial \mathcal{L}}{\partial z_{j_2;\tilde{\alpha}_2}} \frac{\partial \mathcal{L}}{\partial z_{j_3;\tilde{\alpha}_3}} \\ & + \dots \quad \text{ddNTK} \quad ddH(t) \end{aligned}$$

One Aside on Gradient Descent

$$\theta_\mu(t+1) = \theta_\mu(t) - \eta \sum_\nu \lambda_{\mu\nu} \left(\sum \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right)$$

Taylor expansion:

$$\begin{aligned}
 z_{i;\delta}(t+1) = & z_{i;\delta}(t) \quad \text{Neural Tangent Kernel (NTK)} \quad H(t) \\
 & - \eta \sum_{j,\tilde{\alpha}} \left(\sum_{\mu,\nu} \lambda_{\mu\nu} \frac{dz_{i;\delta}}{d\theta_\mu} \frac{dz_{j;\tilde{\alpha}}}{d\theta_\nu} \right) \frac{\partial \mathcal{L}}{\partial z_{j;\tilde{\alpha}}} \quad \text{differential of NTK (dNTK)} \quad dH(t) \\
 O(1/n) \left(& + \frac{\eta^2}{2} \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2} \left(\sum_{\mu_1,\nu_1,\mu_2,\nu_2} \lambda_{\mu_1\nu_1} \lambda_{\mu_2\nu_2} \frac{d^2 z_{i;\delta}}{d\theta_{\mu_1} d\theta_{\mu_2}} \frac{dz_{j_1;\tilde{\alpha}_1}}{d\theta_{\nu_1}} \frac{dz_{j_2;\tilde{\alpha}_2}}{d\theta_{\nu_2}} \right) \frac{\partial \mathcal{L}}{\partial z_{j_1;\tilde{\alpha}_1}} \frac{\partial \mathcal{L}}{\partial z_{j_2;\tilde{\alpha}_2}} \right. \\
 & - \frac{\eta^3}{6} \sum \left(\sum \lambda_{\mu_1\nu_1} \lambda_{\mu_2\nu_2} \lambda_{\mu_3\nu_3} \frac{d^3 z_{i;\delta}}{d\theta_{\mu_1} d\theta_{\mu_2} d\theta_{\mu_3}} \frac{dz_{j_1;\tilde{\alpha}_1}}{d\theta_{\nu_1}} \frac{dz_{j_2;\tilde{\alpha}_2}}{d\theta_{\nu_2}} \frac{dz_{j_3;\tilde{\alpha}_3}}{d\theta_{\nu_3}} \right) \frac{\partial \mathcal{L}}{\partial z_{j_1;\tilde{\alpha}_1}} \frac{\partial \mathcal{L}}{\partial z_{j_2;\tilde{\alpha}_2}} \frac{\partial \mathcal{L}}{\partial z_{j_3;\tilde{\alpha}_3}} \\
 & \left. + \dots \right) \quad \text{ddNTK} \quad ddH(t) \\
 O(1/n^2) \left(& \right.
 \end{aligned}$$

Neural Tangent Kernel (NTK)

$$\hat{H}_{i_1 i_2; \delta_1 \delta_2}^{(\ell)} \equiv \sum_{\mu, \nu} \lambda_{\mu \nu} \frac{d\hat{z}_{i_1; \delta_1}^{(\ell)}}{d\theta_{\mu}} \frac{d\hat{z}_{i_2; \delta_2}^{(\ell)}}{d\theta_{\nu}} \quad \{\theta_{\mu}\} = \{b_i^{(\ell)}, W_{ij}^{(\ell)}\}$$

$$\hat{H}_{i_1 i_2; \delta_1 \delta_2} = \hat{H}_{i_1 i_2; \delta_1 \delta_2}^{(L)}$$

Neural Tangent Kernel (NTK)

$$\hat{H}_{i_1 i_2; \delta_1 \delta_2}^{(\ell)} \equiv \sum_{\mu, \nu} \lambda_{\mu \nu} \frac{d\hat{z}_{i_1; \delta_1}^{(\ell)}}{d\theta_{\mu}} \frac{d\hat{z}_{i_2; \delta_2}^{(\ell)}}{d\theta_{\nu}} \quad \{\theta_{\mu}\} = \{b_i^{(\ell)}, W_{ij}^{(\ell)}\}$$

$$\hat{H}_{i_1 i_2; \delta_1 \delta_2} = \hat{H}_{i_1 i_2; \delta_1 \delta_2}^{(L)}$$

Diagonal, group-by-group, learning rate:

$$\lambda_{b_{i_1}^{(\ell)} b_{i_2}^{(\ell)}} = \delta_{i_1 i_2} \lambda_b^{(\ell)}, \quad \lambda_{W_{i_1 j_1}^{(\ell)} W_{i_2 j_2}^{(\ell)}} = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{\lambda_W^{(\ell)}}{n_{\ell-1}}$$

[Cf. Andrea's "S" matrix]

Neural Tangent Kernel (NTK)

$$\hat{H}_{i_1 i_2; \delta_1 \delta_2}^{(\ell)} \equiv \sum_{\mu, \nu} \lambda_{\mu \nu} \frac{d\hat{z}_{i_1; \delta_1}^{(\ell)}}{d\theta_{\mu}} \frac{d\hat{z}_{i_2; \delta_2}^{(\ell)}}{d\theta_{\nu}} \quad \{\theta_{\mu}\} = \{b_i^{(\ell)}, W_{ij}^{(\ell)}\}$$

$$\hat{H}_{i_1 i_2; \delta_1 \delta_2} = \hat{H}_{i_1 i_2; \delta_1 \delta_2}^{(L)}$$

Diagonal, group-by-group, learning rate:

$$\lambda_{b_{i_1}^{(\ell)} b_{i_2}^{(\ell)}} = \delta_{i_1 i_2} \lambda_b^{(\ell)}, \quad \lambda_{W_{i_1 j_1}^{(\ell)} W_{i_2 j_2}^{(\ell)}} = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{\lambda_W^{(\ell)}}{n_{\ell-1}}$$

good wide limit

Two Pedagogical Simplifications

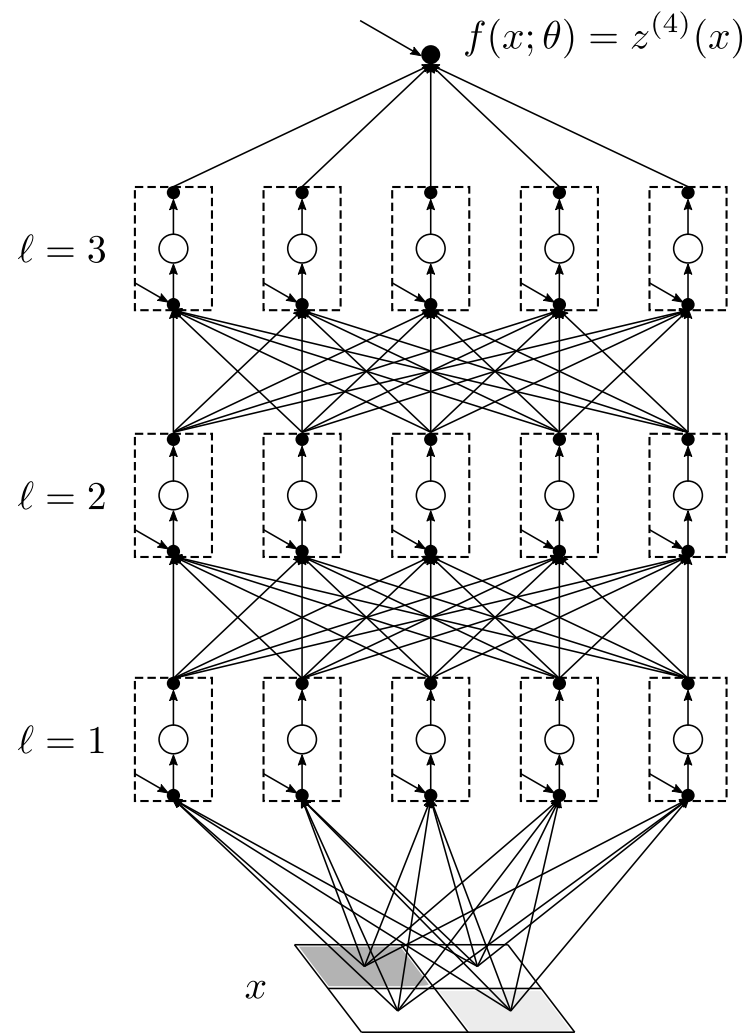
[See “PDLT” (arXiv:2106.10165) for more general cases.]

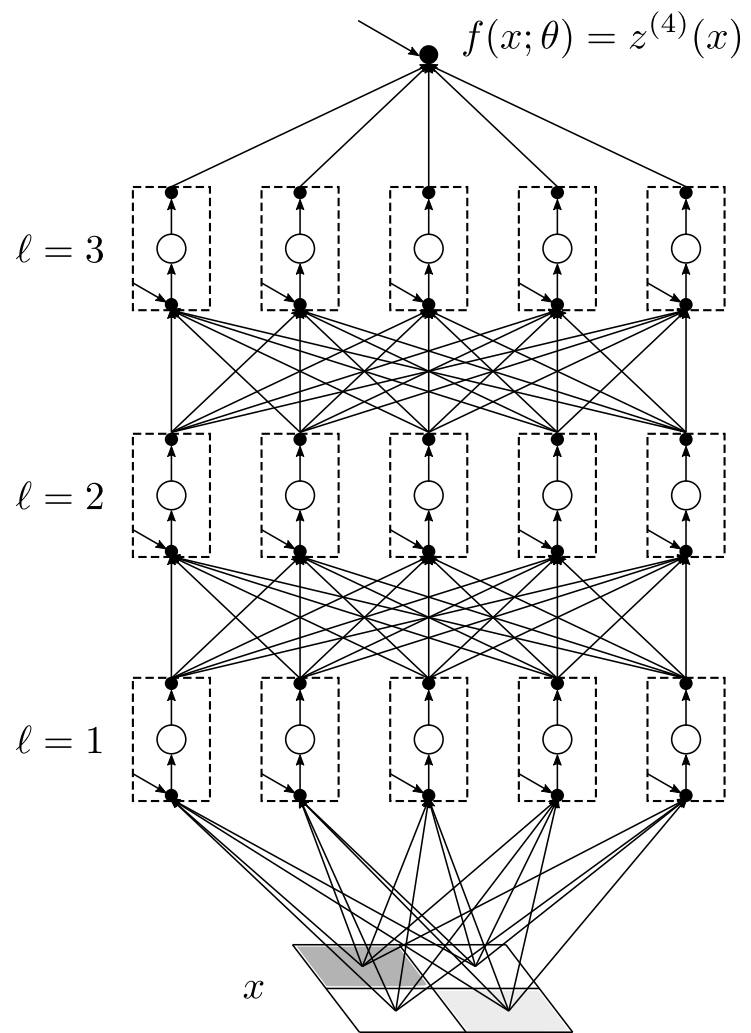
1. Single input; drop sample indices

$$x_{j;\delta} \rightarrow x_j, \quad \hat{z}_{j;\delta}^{(\ell)} \rightarrow \hat{z}_j^{(\ell)}, \quad \hat{H}_{i_1 i_2; \delta_1 \delta_2}^{(\ell)} \rightarrow \hat{H}_{i_1 i_2}^{(\ell)}$$

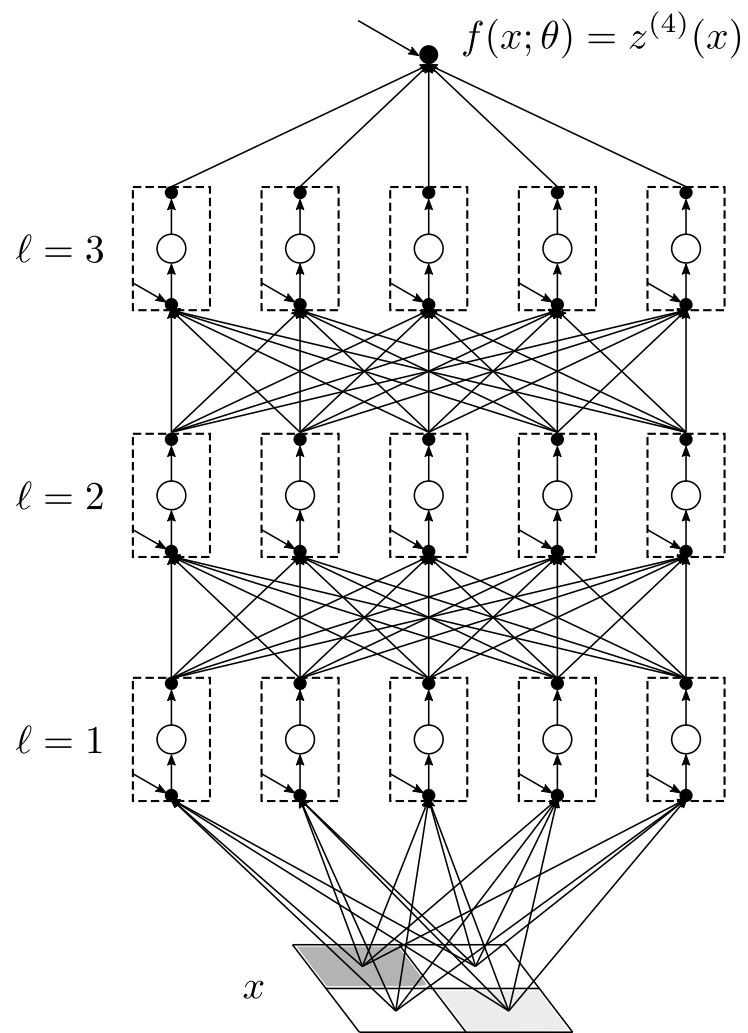
2. Layer-independent hyperparameters; drop layer indices from them

$$C_b^{(\ell)} = C_b, \quad C_W^{(\ell)} = C_W, \quad \lambda_b^{(\ell)} = \lambda_b, \quad \lambda_W^{(\ell)} = \lambda_W$$



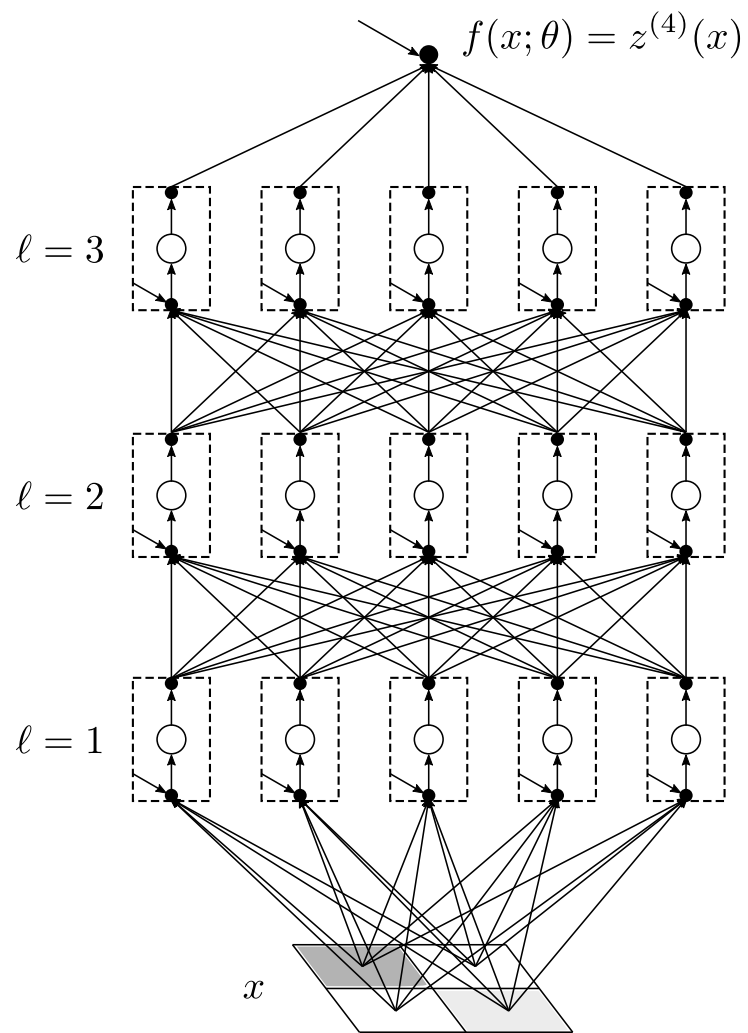


$$p(\hat{z}^{(1)}, \hat{H}^{(1)}, \widehat{dH}^{(1)}, \dots)$$



$$p(\hat{z}^{(2)}, \hat{H}^{(2)}, \widehat{dH}^{(2)}, \dots)$$

$$p(\hat{z}^{(1)}, \hat{H}^{(1)}, \widehat{dH}^{(1)}, \dots)$$



$$p(\hat{z}^{(4)}, \hat{H}^{(4)}, \widehat{\mathrm{d}H}^{(4)}, \dots)$$

$$p(\hat{z}^{(3)}, \hat{H}^{(3)}, \widehat{\mathrm{d}H}^{(3)}, \dots)$$

$$p(\hat{z}^{(2)}, \hat{H}^{(2)}, \widehat{\mathrm{d}H}^{(2)}, \dots)$$

$$p(\hat{z}^{(1)}, \hat{H}^{(1)}, \widehat{\mathrm{d}H}^{(1)}, \dots)$$

2. One-Layer Neural Networks

$$p(\widehat{z}^{(1)}, \widehat{H}^{(1)}, \widehat{dH}^{(1)}, \dots)$$

Statistics of $\hat{z}_i^{(1)} = b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j$

$$p(\hat{z}^{(1)})$$

$$\mathbb{E}[\hat{z}_i^{(1)}], \mathbb{E}[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)}], \mathbb{E}[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)} \hat{z}_{i_3}^{(1)}], \mathbb{E}[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)} \hat{z}_{i_3}^{(1)} \hat{z}_{i_4}^{(1)}], \dots$$

Statistics of $\hat{z}_i^{(1)} = b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j$

$$p(\hat{z}^{(1)})$$

$$\cancel{\mathbb{E}[\hat{z}_i^{(1)}]}, \mathbb{E}[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)}], \cancel{\mathbb{E}[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)} \hat{z}_{i_3}^{(1)}]}, \mathbb{E}[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)} \hat{z}_{i_3}^{(1)} \hat{z}_{i_4}^{(1)}], \dots$$

Statistics of $\widehat{z}_i^{(1)} = b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j$

$$\mathbb{E} \left[\widehat{z}_{i_1}^{(1)} \widehat{z}_{i_2}^{(1)} \right] = \mathbb{E} \left[\left(b_{i_1}^{(1)} + \sum_{j_1=1}^{n_0} W_{i_1 j_1}^{(1)} x_{j_1} \right) \left(b_{i_2}^{(1)} + \sum_{j_2=1}^{n_0} W_{i_2 j_2}^{(1)} x_{j_2} \right) \right]$$

$$\mathbb{E} \left[b_{i_1}^{(1)} b_{i_2}^{(1)} \right] = \delta_{i_1 i_2} C_b, \quad \mathbb{E} \left[W_{i_1 j_1}^{(1)} W_{i_2 j_2}^{(1)} \right] = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{C_W}{n_0}$$

Statistics of $\hat{z}_i^{(1)} = b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j$

“Wick contraction”

$$\mathbb{E} \left[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)} \right] = \mathbb{E} \left[\left(b_{i_1}^{(1)} + \sum_{j_1=1}^{n_0} W_{i_1 j_1}^{(1)} x_{j_1} \right) \left(b_{i_2}^{(1)} + \sum_{j_2=1}^{n_0} W_{i_2 j_2}^{(1)} x_{j_2} \right) \right]$$

$$= \underbrace{C_b \delta_{i_1 i_2}}_{\text{Wick contraction}} + \sum_{j_1, j_2=1}^{n_0} \frac{C_W}{n_0} \delta_{i_1 i_2} \delta_{j_1 j_2} x_{j_1} x_{j_2}$$

$$\mathbb{E} \left[b_{i_1}^{(1)} b_{i_2}^{(1)} \right] = \delta_{i_1 i_2} C_b, \quad \mathbb{E} \left[W_{i_1 j_1}^{(1)} W_{i_2 j_2}^{(1)} \right] = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{C_W}{n_0}$$

Statistics of $\hat{z}_i^{(1)} = b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j$

$$\begin{aligned} \mathbb{E} \left[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)} \right] &= \mathbb{E} \left[\left(b_{i_1}^{(1)} + \sum_{j_1=1}^{n_0} W_{i_1 j_1}^{(1)} x_{j_1} \right) \left(b_{i_2}^{(1)} + \sum_{j_2=1}^{n_0} W_{i_2 j_2}^{(1)} x_{j_2} \right) \right] \\ &= C_b \delta_{i_1 i_2} + \sum_{j_1, j_2=1}^{n_0} \frac{C_W}{n_0} \delta_{i_1 i_2} \delta_{j_1 j_2} x_{j_1} x_{j_2} \end{aligned}$$

$$\mathbb{E} \left[b_{i_1}^{(1)} b_{i_2}^{(1)} \right] = \delta_{i_1 i_2} C_b, \quad \mathbb{E} \left[W_{i_1 j_1}^{(1)} W_{i_2 j_2}^{(1)} \right] = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{C_W}{n_0}$$

Statistics of $\hat{z}_i^{(1)} = b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j$

$$\begin{aligned}\mathbb{E} \left[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)} \right] &= \mathbb{E} \left[\left(b_{i_1}^{(1)} + \sum_{j_1=1}^{n_0} W_{i_1 j_1}^{(1)} x_{j_1} \right) \left(b_{i_2}^{(1)} + \sum_{j_2=1}^{n_0} W_{i_2 j_2}^{(1)} x_{j_2} \right) \right] \\ &= C_b \delta_{i_1 i_2} + \sum_{j_1, j_2=1}^{n_0} \frac{C_W}{n_0} \delta_{i_1 i_2} \delta_{j_1 j_2} x_{j_1} x_{j_2} \\ &= \delta_{i_1 i_2} \left[C_b + C_W \left(\frac{1}{n_0} \sum_{j=1}^{n_0} x_j^2 \right) \right] \equiv \delta_{i_1 i_2} G^{(1)}\end{aligned}$$

Statistics of $\hat{z}_i^{(1)} = b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j$

$$\begin{aligned}
& \mathbb{E} \left[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)} \hat{z}_{i_3}^{(1)} \hat{z}_{i_4}^{(1)} \right] \\
&= \mathbb{E} \left[\left(b_{i_1}^{(1)} + \sum_{j_1=1}^{n_0} W_{i_1 j_1}^{(1)} x_{j_1} \right) \left(b_{i_2}^{(1)} + \sum_{j_2=1}^{n_0} W_{i_2 j_2}^{(1)} x_{j_2} \right) \left(b_{i_3}^{(1)} + \sum_{j_3=1}^{n_0} W_{i_3 j_3}^{(1)} x_{j_3} \right) \left(b_{i_4}^{(1)} + \sum_{j_4=1}^{n_0} W_{i_4 j_4}^{(1)} x_{j_4} \right) \right] \\
&= (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3}) \\
&\quad \times \left(C_b^2 + 2C_b C_W \frac{1}{n_0} \sum_{j=1}^{n_0} x_j^2 + C_W^2 \frac{1}{n_0^2} \sum_{j_1, j_2=0}^{n_0} x_{j_1}^2 x_{j_2}^2 \right) \\
&= \left(G^{(1)} \right)^2 (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3})
\end{aligned}$$

Statistics of $\hat{z}_i^{(1)} = b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j$

$$\begin{aligned}
 & \mathbb{E} \left[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)} \hat{z}_{i_3}^{(1)} \hat{z}_{i_4}^{(1)} \right] \\
 &= \mathbb{E} \left[\left(b_{i_1}^{(1)} + \sum_{j_1=1}^{n_0} W_{i_1 j_1}^{(1)} x_{j_1} \right) \left(b_{i_2}^{(1)} + \sum_{j_2=1}^{n_0} W_{i_2 j_2}^{(1)} x_{j_2} \right) \left(b_{i_3}^{(1)} + \sum_{j_3=1}^{n_0} W_{i_3 j_3}^{(1)} x_{j_3} \right) \left(b_{i_4}^{(1)} + \sum_{j_4=1}^{n_0} W_{i_4 j_4}^{(1)} x_{j_4} \right) \right] \\
 &= (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3}) \\
 &\quad \times \left(\underline{C_b^2} + 2C_b C_W \frac{1}{n_0} \sum_{j=1}^{n_0} x_j^2 + C_W^2 \frac{1}{n_0^2} \sum_{j_1, j_2=0}^{n_0} x_{j_1}^2 x_{j_2}^2 \right) \\
 &= \left(G^{(1)} \right)^2 (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3})
 \end{aligned}$$

$$\mathbb{E} \left[b_{i_1}^{(1)} b_{i_2}^{(1)} \right] = \delta_{i_1 i_2} C_b, \quad \mathbb{E} \left[W_{i_1 j_1}^{(1)} W_{i_2 j_2}^{(1)} \right] = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{C_W}{n_0}$$

Statistics of $\hat{z}_i^{(1)} = b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j$

$$\begin{aligned}
 & \mathbb{E} \left[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)} \hat{z}_{i_3}^{(1)} \hat{z}_{i_4}^{(1)} \right] \\
 &= \mathbb{E} \left[\left(b_{i_1}^{(1)} + \sum_{j_1=1}^{n_0} W_{i_1 j_1}^{(1)} x_{j_1} \right) \left(b_{i_2}^{(1)} + \sum_{j_2=1}^{n_0} W_{i_2 j_2}^{(1)} x_{j_2} \right) \left(b_{i_3}^{(1)} + \sum_{j_3=1}^{n_0} W_{i_3 j_3}^{(1)} x_{j_3} \right) \left(b_{i_4}^{(1)} + \sum_{j_4=1}^{n_0} W_{i_4 j_4}^{(1)} x_{j_4} \right) \right] \\
 &= (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3}) \\
 &\quad \times \left(\underline{C_b^2} + 2C_b C_W \frac{1}{n_0} \sum_{j=1}^{n_0} x_j^2 + C_W^2 \frac{1}{n_0^2} \sum_{j_1, j_2=0}^{n_0} x_{j_1}^2 x_{j_2}^2 \right) \\
 &= \left(G^{(1)} \right)^2 (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3})
 \end{aligned}$$

$$\mathbb{E} \left[b_{i_1}^{(1)} b_{i_2}^{(1)} \right] = \delta_{i_1 i_2} C_b, \quad \mathbb{E} \left[W_{i_1 j_1}^{(1)} W_{i_2 j_2}^{(1)} \right] = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{C_W}{n_0}$$

Statistics of $\hat{z}_i^{(1)} = b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j$

$$\begin{aligned}
 & \mathbb{E} \left[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)} \hat{z}_{i_3}^{(1)} \hat{z}_{i_4}^{(1)} \right] \\
 &= \mathbb{E} \left[\left(b_{i_1}^{(1)} + \sum_{j_1=1}^{n_0} W_{i_1 j_1}^{(1)} x_{j_1} \right) \left(b_{i_2}^{(1)} + \sum_{j_2=1}^{n_0} W_{i_2 j_2}^{(1)} x_{j_2} \right) \left(b_{i_3}^{(1)} + \sum_{j_3=1}^{n_0} W_{i_3 j_3}^{(1)} x_{j_3} \right) \left(b_{i_4}^{(1)} + \sum_{j_4=1}^{n_0} W_{i_4 j_4}^{(1)} x_{j_4} \right) \right] \\
 &= (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} - \delta_{i_1 i_4} \delta_{i_2 i_3}) \\
 &\quad \times \left(\underline{C_b^2} + 2C_b C_W \frac{1}{n_0} \sum_{j=1}^{n_0} x_j^2 + C_W^2 \frac{1}{n_0^2} \sum_{j_1, j_2=0}^{n_0} x_{j_1}^2 x_{j_2}^2 \right) \\
 &= \left(G^{(1)} \right)^2 (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3})
 \end{aligned}$$

$$\mathbb{E} \left[b_{i_1}^{(1)} b_{i_2}^{(1)} \right] = \delta_{i_1 i_2} C_b, \quad \mathbb{E} \left[W_{i_1 j_1}^{(1)} W_{i_2 j_2}^{(1)} \right] = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{C_W}{n_0}$$

Statistics of $\hat{z}_i^{(1)} = b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j$

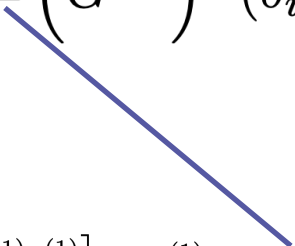
$$\begin{aligned}
 & \mathbb{E} \left[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)} \hat{z}_{i_3}^{(1)} \hat{z}_{i_4}^{(1)} \right] \\
 &= \mathbb{E} \left[\left(b_{i_1}^{(1)} + \sum_{j_1=1}^{n_0} W_{i_1 j_1}^{(1)} x_{j_1} \right) \left(b_{i_2}^{(1)} + \sum_{j_2=1}^{n_0} W_{i_2 j_2}^{(1)} x_{j_2} \right) \left(b_{i_3}^{(1)} + \sum_{j_3=1}^{n_0} W_{i_3 j_3}^{(1)} x_{j_3} \right) \left(b_{i_4}^{(1)} + \sum_{j_4=1}^{n_0} W_{i_4 j_4}^{(1)} x_{j_4} \right) \right] \\
 &= (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3}) \\
 &\quad \times \left(C_b^2 + 2C_b C_W \frac{1}{n_0} \sum_{j=1}^{n_0} x_j^2 + C_W^2 \frac{1}{n_0^2} \sum_{j_1, j_2=0}^{n_0} x_{j_1}^2 x_{j_2}^2 \right) \\
 &= \left(G^{(1)} \right)^2 (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3})
 \end{aligned}$$

$bbWW$
 $WWWW$

$$\mathbb{E} \left[b_{i_1}^{(1)} b_{i_2}^{(1)} \right] = \delta_{i_1 i_2} C_b, \quad \mathbb{E} \left[W_{i_1 j_1}^{(1)} W_{i_2 j_2}^{(1)} \right] = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{C_W}{n_0}$$

Statistics of $\hat{z}_i^{(1)} = b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j$

$$\begin{aligned}
 & \mathbb{E} \left[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)} \hat{z}_{i_3}^{(1)} \hat{z}_{i_4}^{(1)} \right] \\
 &= \mathbb{E} \left[\left(b_{i_1}^{(1)} + \sum_{j_1=1}^{n_0} W_{i_1 j_1}^{(1)} x_{j_1} \right) \left(b_{i_2}^{(1)} + \sum_{j_2=1}^{n_0} W_{i_2 j_2}^{(1)} x_{j_2} \right) \left(b_{i_3}^{(1)} + \sum_{j_3=1}^{n_0} W_{i_3 j_3}^{(1)} x_{j_3} \right) \left(b_{i_4}^{(1)} + \sum_{j_4=1}^{n_0} W_{i_4 j_4}^{(1)} x_{j_4} \right) \right] \\
 &= (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3}) \\
 &\quad \times \left(C_b^2 + 2C_b C_W \frac{1}{n_0} \sum_{j=1}^{n_0} x_j^2 + C_W^2 \frac{1}{n_0^2} \sum_{j_1, j_2=0}^{n_0} x_{j_1}^2 x_{j_2}^2 \right) \\
 &= \left(G^{(1)} \right)^2 (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3})
 \end{aligned}$$



$$\mathbb{E} \left[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)} \right] = G^{(1)} \delta_{i_1 i_2} = \left[C_b + C_W \left(\frac{1}{n_0} \sum_j x_j^2 \right) \right] \delta_{i_1 i_2}$$

Statistics of $\hat{z}_i^{(1)} = b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j$

$$\mathbb{E} \left[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)} \right] = G^{(1)} \delta_{i_1 i_2}$$

$$\mathbb{E} \left[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)} \hat{z}_{i_3}^{(1)} \hat{z}_{i_4}^{(1)} \right] = \left(G^{(1)} \right)^2 \left(\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3} \right)$$

...

Statistics of $\hat{z}_i^{(1)} = b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j$

$$\mathbb{E} \left[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)} \right] = G^{(1)} \delta_{i_1 i_2}$$

$$\mathbb{E} \left[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)} \hat{z}_{i_3}^{(1)} \hat{z}_{i_4}^{(1)} \right] = \left(G^{(1)} \right)^2 \left(\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3} \right)$$

...

$$p\left(\hat{z}^{(1)}\right) \propto \exp \left[-\frac{1}{2G^{(1)}} \sum_{i=1}^{n_1} \left(\hat{z}_i^{(1)} \right)^2 \right] = \prod_{i=1}^{n_1} \left\{ \exp \left[-\frac{1}{2G^{(1)}} \left(\hat{z}_i^{(1)} \right)^2 \right] \right\}$$

Statistics of $\hat{z}_i^{(1)} = b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j$

$$p(\hat{z}^{(1)}) \propto \exp \left[-\frac{1}{2G^{(1)}} \sum_{i=1}^{n_1} \left(\hat{z}_i^{(1)} \right)^2 \right] = \prod_{i=1}^{n_1} \left\{ \exp \left[-\frac{1}{2G^{(1)}} \left(\hat{z}_i^{(1)} \right)^2 \right] \right\}$$

- Neurons don't talk to each other; they are statistically independent.
- We marginalized over/integrated out $b_i^{(1)}$ and $W_{ij}^{(1)}$.
- Two interpretations:
 - (i) outputs of one-layer networks; or
 - (ii) preactivations in the first layer of deeper networks.

Statistics of $\hat{H}_{i_1 i_2}^{(1)}$

$$\hat{H}_{i_1 i_2}^{(1)} \equiv \sum_{\mu, \nu} \lambda_{\mu \nu} \frac{d\hat{z}_{i_1}^{(1)}}{d\theta_{\mu}} \frac{d\hat{z}_{i_2}^{(1)}}{d\theta_{\nu}}$$

Statistics of $\hat{H}_{i_1 i_2}^{(1)}$

$$\begin{aligned}\hat{H}_{i_1 i_2}^{(1)} &\equiv \sum_{\mu, \nu} \lambda_{\mu \nu} \frac{d\hat{z}_{i_1}^{(1)}}{d\theta_{\mu}} \frac{d\hat{z}_{i_2}^{(1)}}{d\theta_{\nu}} \\ &= \lambda_b \sum_{j=1}^{n_1} \frac{d\hat{z}_{i_1}^{(1)}}{db_j^{(1)}} \frac{d\hat{z}_{i_2}^{(1)}}{db_j^{(1)}} + \frac{\lambda_W}{n_0} \sum_{j=1}^{n_1} \sum_{k=1}^{n_0} \frac{d\hat{z}_{i_1}^{(1)}}{dW_{jk}^{(1)}} \frac{d\hat{z}_{i_2}^{(1)}}{dW_{jk}^{(1)}}\end{aligned}$$

$$\lambda_{b_{i_1}^{(1)} b_{i_2}^{(1)}} = \delta_{i_1 i_2} \lambda_b, \quad \lambda_{W_{i_1 j_1}^{(1)} W_{i_2 j_2}^{(1)}} = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{\lambda_W}{n_0}$$

Statistics of $\hat{H}_{i_1 i_2}^{(1)}$

$$\begin{aligned}
 \hat{H}_{i_1 i_2}^{(1)} &\equiv \sum_{\mu, \nu} \lambda_{\mu \nu} \frac{d\hat{z}_{i_1}^{(1)}}{d\theta_{\mu}} \frac{d\hat{z}_{i_2}^{(1)}}{d\theta_{\nu}} \\
 &= \lambda_b \sum_{j=1}^{n_1} \frac{d\hat{z}_{i_1}^{(1)}}{db_j^{(1)}} \frac{d\hat{z}_{i_2}^{(1)}}{db_j^{(1)}} + \frac{\lambda_W}{n_0} \sum_{j=1}^{n_1} \sum_{k=1}^{n_0} \frac{d\hat{z}_{i_1}^{(1)}}{dW_{jk}^{(1)}} \frac{d\hat{z}_{i_2}^{(1)}}{dW_{jk}^{(1)}} \\
 &= \lambda_b \sum_{j=1}^{n_1} \delta_{i_1 j} \delta_{i_2 j} + \frac{\lambda_W}{n_0} \sum_{j=1}^{n_1} \sum_{k=1}^{n_0} \delta_{i_1 j} x_k \delta_{i_2 j} x_k
 \end{aligned}$$

$$\hat{z}_i^{(1)} = b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j$$

Statistics of $\hat{H}_{i_1 i_2}^{(1)}$

$$\begin{aligned}
 \hat{H}_{i_1 i_2}^{(1)} &\equiv \sum_{\mu, \nu} \lambda_{\mu \nu} \frac{d\hat{z}_{i_1}^{(1)}}{d\theta_{\mu}} \frac{d\hat{z}_{i_2}^{(1)}}{d\theta_{\nu}} \\
 &= \lambda_b \sum_{j=1}^{n_1} \frac{d\hat{z}_{i_1}^{(1)}}{db_j^{(1)}} \frac{d\hat{z}_{i_2}^{(1)}}{db_j^{(1)}} + \frac{\lambda_W}{n_0} \sum_{j=1}^{n_1} \sum_{k=1}^{n_0} \frac{d\hat{z}_{i_1}^{(1)}}{dW_{jk}^{(1)}} \frac{d\hat{z}_{i_2}^{(1)}}{dW_{jk}^{(1)}} \\
 &= \lambda_b \sum_{j=1}^{n_1} \delta_{i_1 j} \delta_{i_2 j} + \frac{\lambda_W}{n_0} \sum_{j=1}^{n_1} \sum_{k=1}^{n_0} \delta_{i_1 j} x_k \delta_{i_2 j} x_k \\
 &= \lambda_b \delta_{i_1 i_2} + \frac{\lambda_W}{n_0} \delta_{i_1 i_2} \sum_{k=1}^{n_0} x_k x_k \\
 &= \delta_{i_1 i_2} \left[\lambda_b + \lambda_W \left(\frac{1}{n_0} \sum_{j=1}^{n_0} x_j^2 \right) \right] \equiv \delta_{i_1 i_2} H^{(1)}
 \end{aligned}$$

Statistics of $\hat{H}_{i_1 i_2}^{(1)}$

$$\hat{H}_{i_1 i_2}^{(1)} = \delta_{i_1 i_2} H^{(1)} = \delta_{i_1 i_2} \left[\lambda_b + \lambda_W \left(\frac{1}{n_0} \sum_{j=1}^{n_0} x_j^2 \right) \right]$$

- “Deterministic”: it doesn't depend on any particular initialization; you always get the same number.
- “Frozen”: it cannot evolve during training; no representation learning.

Statistics of $\widehat{\mathrm{d}H}_{i_0 i_1 i_2}^{(1)}$

$$\widehat{\mathrm{d}H}_{i_0 i_1 i_2}^{(1)} \equiv \sum_{\mu_1, \nu_1, \mu_2, \nu_2} \lambda_{\mu_1 \nu_1} \lambda_{\mu_2 \nu_2} \frac{d^2 \widehat{z}_{i_0}^{(1)}}{d\theta_{\mu_1} d\theta_{\mu_2}} \frac{d\widehat{z}_{i_1}^{(1)}}{d\theta_{\nu_1}} \frac{d\widehat{z}_{i_2}^{(1)}}{d\theta_{\nu_2}}$$

Statistics of $\widehat{dH}_{i_0 i_1 i_2}^{(1)}$

$$\widehat{dH}_{i_0 i_1 i_2}^{(1)} \equiv \sum_{\mu_1, \nu_1, \mu_2, \nu_2} \lambda_{\mu_1 \nu_1} \lambda_{\mu_2 \nu_2} \frac{d^2 \widehat{z}_{i_0}^{(1)}}{d\theta_{\mu_1} d\theta_{\mu_2}} \frac{d\widehat{z}_{i_1}^{(1)}}{d\theta_{\nu_1}} \frac{d\widehat{z}_{i_2}^{(1)}}{d\theta_{\nu_2}} = 0$$

$$\widehat{z}_i^{(1)} = b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j$$

Statistics of $\widehat{dH}_{i_0 i_1 i_2}^{(1)}$

$$\widehat{dH}^{(1)}, \widehat{ddH}^{(1)}, \dots = 0$$

- No representation learning.
- No algorithm dependence.

Statistics of One-Layer Neural Networks

$$p\left(\widehat{z}^{(1)}\right) \propto \exp \left[-\frac{1}{2G^{(1)}} \sum_{i=1}^{n_1} \left(\widehat{z}_i^{(1)}\right)^2\right] = \prod_{i=1}^{n_1} \left\{ \exp \left[-\frac{1}{2G^{(1)}} \left(\widehat{z}_i^{(1)}\right)^2\right] \right\}$$

$$\widehat{H}_{i_1 i_2}^{(1)} = \delta_{i_1 i_2} H^{(1)} = \delta_{i_1 i_2} \left[\lambda_b + \lambda_W \left(\frac{1}{n_0} \sum_{j=1}^{n_0} x_j^2 \right) \right]$$

$$\widehat{dH}^{(1)}, \widehat{ddH}^{(1)}, \dots = 0$$

Statistics of One-Layer Neural Networks

$$p(\hat{z}^{(1)}) \propto \exp \left[-\frac{1}{2G^{(1)}} \sum_{i=1}^{n_1} \left(\hat{z}_i^{(1)} \right)^2 \right] = \prod_{i=1}^{n_1} \left\{ \exp \left[-\frac{1}{2G^{(1)}} \left(\hat{z}_i^{(1)} \right)^2 \right] \right\}$$

$$\hat{H}_{i_1 i_2}^{(1)} = \delta_{i_1 i_2} H^{(1)} = \delta_{i_1 i_2} \left[\lambda_b + \lambda_W \left(\frac{1}{n_0} \sum_{j=1}^{n_0} x_j^2 \right) \right]$$

$$\widehat{dH}^{(1)}, \widehat{ddH}^{(1)}, \dots = 0$$

Linear dynamics:
$$z_{i;\delta}(t+1) = z_{i;\delta}(t) - \eta \sum H_{\delta\tilde{\alpha}} [z_{i;\tilde{\alpha}}(t) - y_{i;\tilde{\alpha}}]$$

Simple solution:
$$z_{i;\delta}^{\star} = \hat{z}_{i;\delta} - \sum H_{\delta\tilde{\alpha}_1}^{(1)} \left(\left(\tilde{H}^{(1)} \right)^{-1} \right)^{\tilde{\alpha}_1 \tilde{\alpha}_2} [\hat{z}_{i;\tilde{\alpha}_2} - y_{i;\tilde{\alpha}_2}]$$

Statistics of One-Layer Neural Networks

$$p(\hat{z}^{(1)}) \propto \exp \left[-\frac{1}{2G^{(1)}} \sum_{i=1}^{n_1} \left(\hat{z}_i^{(1)} \right)^2 \right] = \prod_{i=1}^{n_1} \left\{ \exp \left[-\frac{1}{2G^{(1)}} \left(\hat{z}_i^{(1)} \right)^2 \right] \right\}$$

$$\hat{H}_{i_1 i_2}^{(1)} = \delta_{i_1 i_2} H^{(1)} = \delta_{i_1 i_2} \left[\lambda_b + \lambda_W \left(\frac{1}{n_0} \sum_{j=1}^{n_0} x_j^2 \right) \right]$$

$$\widehat{dH}^{(1)}, \widehat{ddH}^{(1)}, \dots = 0$$

$$z_{i;\delta}^* = \hat{z}_{i;\delta} - \sum H_{\delta \tilde{\alpha}_1}^{(1)} \left(\left(\tilde{H}^{(1)} \right)^{-1} \right)^{\tilde{\alpha}_1 \tilde{\alpha}_2} \left[\hat{z}_{i;\tilde{\alpha}_2} - y_{i;\tilde{\alpha}_2} \right]$$

$$p(\hat{z}, H) \rightarrow p(z^*)$$

statistics at *initialization*

statistics *after training*

$$\mathbb{E} [z_{i;\delta}^*] = \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2 \in \mathcal{A}} H_{\delta \tilde{\alpha}_1}^{(1)} \left(\left(H^{(1)} \right)^{-1} \right)^{\tilde{\alpha}_1 \tilde{\alpha}_2} y_{i;\tilde{\alpha}_2}$$

Statistics of One-Layer Neural Networks

$$p(\hat{z}^{(1)}) \propto \exp \left[-\frac{1}{2G^{(1)}} \sum_{i=1}^{n_1} \left(\hat{z}_i^{(1)} \right)^2 \right] = \prod_{i=1}^{n_1} \left\{ \exp \left[-\frac{1}{2G^{(1)}} \left(\hat{z}_i^{(1)} \right)^2 \right] \right\}$$

$$\hat{H}_{i_1 i_2}^{(1)} = \delta_{i_1 i_2} H^{(1)} = \delta_{i_1 i_2} \left[\lambda_b + \lambda_W \left(\frac{1}{n_0} \sum_{j=1}^{n_0} x_j^2 \right) \right]$$

$$\widehat{dH}^{(1)}, \widehat{ddH}^{(1)}, \dots = 0$$

- Same trivial statistics for infinite-width neural networks of *any* fixed depth.
- No representation learning, no algorithm dependence; not a good model of deep learning.

We must study deeper networks of finite width!

3. Two-Layer Neural Networks

$$p(\hat{z}^{(2)}, \hat{H}^{(2)}, \widehat{dH}^{(2)}, \dots)$$

Statistics of $\hat{z}_i^{(2)} = b_i^{(2)} + \sum_{j=1}^{n_1} W_{ij}^{(2)} \sigma \left(\hat{z}_j^{(1)} \right)$

$$\mathbb{E} \left[\hat{z}_{i_1}^{(2)} \hat{z}_{i_2}^{(2)} \right] = \mathbb{E} \left[\left(b_{i_1}^{(2)} + \sum_{j_1=1}^{n_1} W_{i_1 j_1}^{(2)} \sigma \left(\hat{z}_{j_1}^{(1)} \right) \right) \left(b_{i_2}^{(2)} + \sum_{j_2=1}^{n_1} W_{i_2 j_2}^{(2)} \sigma \left(\hat{z}_{j_2}^{(1)} \right) \right) \right]$$

$$\mathbb{E} \left[b_{i_1}^{(2)} b_{i_2}^{(2)} \right] = \delta_{i_1 i_2} C_b, \quad \mathbb{E} \left[W_{i_1 j_1}^{(2)} W_{i_2 j_2}^{(2)} \right] = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{C_W}{n_1}$$

Statistics of $\hat{z}_i^{(2)} = b_i^{(2)} + \sum_{j=1}^{n_1} W_{ij}^{(2)} \sigma \left(\hat{z}_j^{(1)} \right)$

$$\mathbb{E} \left[\hat{z}_{i_1}^{(2)} \hat{z}_{i_2}^{(2)} \right] = \mathbb{E} \left[\left(b_{i_1}^{(2)} + \sum_{j_1=1}^{n_1} W_{i_1 j_1}^{(2)} \sigma \left(\hat{z}_{j_1}^{(1)} \right) \right) \left(b_{i_2}^{(2)} + \sum_{j_2=1}^{n_1} W_{i_2 j_2}^{(2)} \sigma \left(\hat{z}_{j_2}^{(1)} \right) \right) \right]$$

Wick $= C_b \delta_{i_1 i_2} + \sum_{j_1, j_2=1}^{n_1} \frac{C_W}{n_1} \delta_{i_1 i_2} \delta_{j_1 j_2} \mathbb{E} \left[\sigma \left(\hat{z}_{j_1}^{(1)} \right) \sigma \left(\hat{z}_{j_2}^{(1)} \right) \right]$

arrange $= \delta_{i_1 i_2} \left[C_b + C_W \left(\frac{1}{n_1} \sum_{j=1}^{n_1} \mathbb{E} \left[\sigma \left(\hat{z}_j^{(1)} \right) \sigma \left(\hat{z}_j^{(1)} \right) \right] \right) \right]$

$$\mathbb{E} \left[b_{i_1}^{(2)} b_{i_2}^{(2)} \right] = \delta_{i_1 i_2} C_b, \quad \mathbb{E} \left[W_{i_1 j_1}^{(2)} W_{i_2 j_2}^{(2)} \right] = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{C_W}{n_1}$$

Statistics of $\hat{z}_i^{(2)} = b_i^{(2)} + \sum_{j=1}^{n_1} W_{ij}^{(2)} \sigma \left(\hat{z}_j^{(1)} \right)$

$$\begin{aligned}
 \mathbb{E} \left[\hat{z}_{i_1}^{(2)} \hat{z}_{i_2}^{(2)} \right] &= \mathbb{E} \left[\left(b_{i_1}^{(2)} + \sum_{j_1=1}^{n_1} W_{i_1 j_1}^{(2)} \sigma \left(\hat{z}_{j_1}^{(1)} \right) \right) \left(b_{i_2}^{(2)} + \sum_{j_2=1}^{n_1} W_{i_2 j_2}^{(2)} \sigma \left(\hat{z}_{j_2}^{(1)} \right) \right) \right] \\
 &= C_b \delta_{i_1 i_2} + \sum_{j_1, j_2=1}^{n_1} \frac{C_W}{n_1} \delta_{i_1 i_2} \delta_{j_1 j_2} \mathbb{E} \left[\sigma \left(\hat{z}_{j_1}^{(1)} \right) \sigma \left(\hat{z}_{j_2}^{(1)} \right) \right] \\
 &= \delta_{i_1 i_2} \left[C_b + C_W \left(\frac{1}{n_1} \sum_{j=1}^{n_1} \mathbb{E} \left[\sigma \left(\hat{z}_j^{(1)} \right) \sigma \left(\hat{z}_j^{(1)} \right) \right] \right) \right] \\
 &= \delta_{i_1 i_2} \left[C_b + C_W \langle \sigma(z) \sigma(z) \rangle_{G^{(1)}} \right] \equiv \delta_{i_1 i_2} G^{(2)}
 \end{aligned}$$

$$p(\hat{z}^{(1)}) \propto \exp \left[-\frac{1}{2G^{(1)}} \sum_{i=1}^{n_1} \left(\hat{z}_i^{(1)} \right)^2 \right] = \prod_{i=1}^{n_1} \left\{ \exp \left[-\frac{1}{2G^{(1)}} \left(\hat{z}_i^{(1)} \right)^2 \right] \right\}$$

$$\langle f(z) \rangle_G \equiv \frac{1}{\sqrt{2\pi G}} \int dz f(z) e^{-\frac{z^2}{2G}}$$

Statistics of $\hat{z}_i^{(2)} = b_i^{(2)} + \sum_{j=1}^{n_1} W_{ij}^{(2)} \sigma(\hat{z}_j^{(1)})$

$$\begin{aligned}
 \mathbb{E} \left[\hat{z}_{i_1}^{(2)} \hat{z}_{i_2}^{(2)} \right] &= \mathbb{E} \left[\left(b_{i_1}^{(2)} + \sum_{j_1=1}^{n_1} W_{i_1 j_1}^{(2)} \sigma(\hat{z}_{j_1}^{(1)}) \right) \left(b_{i_2}^{(2)} + \sum_{j_2=1}^{n_1} W_{i_2 j_2}^{(2)} \sigma(\hat{z}_{j_2}^{(1)}) \right) \right] \\
 &= C_b \delta_{i_1 i_2} + \sum_{j_1, j_2=1}^{n_1} \frac{C_W}{n_1} \delta_{i_1 i_2} \delta_{j_1 j_2} \mathbb{E} \left[\sigma(\hat{z}_{j_1}^{(1)}) \sigma(\hat{z}_{j_2}^{(1)}) \right] \\
 &= \delta_{i_1 i_2} \left[C_b + C_W \left(\frac{1}{n_1} \sum_{j=1}^{n_1} \mathbb{E} \left[\sigma(\hat{z}_j^{(1)}) \sigma(\hat{z}_j^{(1)}) \right] \right) \right] \\
 &= \delta_{i_1 i_2} [C_b + C_W \langle \sigma(z) \sigma(z) \rangle_{G^{(1)}}] \equiv \delta_{i_1 i_2} G^{(2)}
 \end{aligned}$$

- Recursive.

- $\mathbb{E} \left[W_{i_1 j_1}^{(2)} W_{i_2 j_2}^{(2)} \right] = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{C_W}{n_1}$ width-scaling was important.

Statistics of $\hat{z}_i^{(2)} = b_i^{(2)} + \sum_{j=1}^{n_1} W_{ij}^{(2)} \sigma \left(\hat{z}_j^{(1)} \right)$

$$\mathbb{E} \left[\hat{z}_{i_1}^{(2)} \hat{z}_{i_2}^{(2)} \hat{z}_{i_3}^{(2)} \hat{z}_{i_4}^{(2)} \right]$$

$$= \mathbb{E} \left[\left(b_{i_1}^{(2)} + \sum_{j_1=1}^{n_1} W_{i_1 j_1}^{(2)} \sigma \left(\hat{z}_{j_1}^{(1)} \right) \right) \cdots \left(b_{i_4}^{(2)} + \sum_{j_4=1}^{n_1} W_{i_4 j_4}^{(2)} \sigma \left(\hat{z}_{j_4}^{(1)} \right) \right) \right]$$

Wick = $(\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3})$

$$\times \left\{ C_b^2 + 2C_b C_W \frac{1}{n_1} \sum_{j=1}^{n_1} \mathbb{E} \left[\sigma \left(\hat{z}_j^{(1)} \right) \sigma \left(\hat{z}_j^{(1)} \right) \right] + C_W^2 \frac{1}{n_1^2} \sum_{j_1, j_2=1}^{n_1} \mathbb{E} \left[\sigma \left(\hat{z}_{j_1}^{(1)} \right) \sigma \left(\hat{z}_{j_1}^{(1)} \right) \sigma \left(\hat{z}_{j_2}^{(1)} \right) \sigma \left(\hat{z}_{j_2}^{(1)} \right) \right] \right\}$$

$$\mathbb{E} \left[b_{i_1}^{(2)} b_{i_2}^{(2)} \right] = \delta_{i_1 i_2} C_b, \quad \mathbb{E} \left[W_{i_1 j_1}^{(2)} W_{i_2 j_2}^{(2)} \right] = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{C_W}{n_1}$$

Statistics of $\hat{z}_i^{(2)} = b_i^{(2)} + \sum_{j=1}^{n_1} W_{ij}^{(2)} \sigma \left(\hat{z}_j^{(1)} \right)$

$$\mathbb{E} \left[\hat{z}_{i_1}^{(2)} \hat{z}_{i_2}^{(2)} \hat{z}_{i_3}^{(2)} \hat{z}_{i_4}^{(2)} \right]$$

$$= \mathbb{E} \left[\left(b_{i_1}^{(2)} + \sum_{j_1=1}^{n_1} W_{i_1 j_1}^{(2)} \sigma \left(\hat{z}_{j_1}^{(1)} \right) \right) \cdots \left(b_{i_4}^{(2)} + \sum_{j_4=1}^{n_1} W_{i_4 j_4}^{(2)} \sigma \left(\hat{z}_{j_4}^{(1)} \right) \right) \right]$$

Wick = $(\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3})$

$$\times \left\{ C_b^2 + 2C_b C_W \frac{1}{n_1} \sum_{j=1}^{n_1} \mathbb{E} \left[\sigma \left(\hat{z}_j^{(1)} \right) \sigma \left(\hat{z}_j^{(1)} \right) \right] + C_W^2 \frac{1}{n_1^2} \left(\sum_{j_1, j_2=1}^{n_1} \mathbb{E} \left[\sigma \left(\hat{z}_{j_1}^{(1)} \right) \sigma \left(\hat{z}_{j_1}^{(1)} \right) \sigma \left(\hat{z}_{j_2}^{(1)} \right) \sigma \left(\hat{z}_{j_2}^{(1)} \right) \right] \right) \right\}$$

$$= (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3})$$

$$\times \left\{ C_b^2 + 2C_b C_W \langle \sigma(z) \sigma(z) \rangle_{G^{(1)}} + C_W^2 \left[\frac{n_1^2 - n_1}{n_1^2} \langle \sigma(z) \sigma(z) \rangle_{G^{(1)}} \langle \sigma(z) \sigma(z) \rangle_{G^{(1)}} + \frac{n_1}{n_1^2} \langle \sigma(z) \sigma(z) \sigma(z) \sigma(z) \rangle_{G^{(1)}} \right] \right\}$$

$j_1 \neq j_2$
 $j_1 = j_2$

Statistics of $\hat{z}_i^{(2)} = b_i^{(2)} + \sum_{j=1}^{n_1} W_{ij}^{(2)} \sigma \left(\hat{z}_j^{(1)} \right)$

$$\mathbb{E} \left[\hat{z}_{i_1}^{(2)} \hat{z}_{i_2}^{(2)} \hat{z}_{i_3}^{(2)} \hat{z}_{i_4}^{(2)} \right]$$

$$= \mathbb{E} \left[\left(b_{i_1}^{(2)} + \sum_{j_1=1}^{n_1} W_{i_1 j_1}^{(2)} \sigma \left(\hat{z}_{j_1}^{(1)} \right) \right) \cdots \left(b_{i_4}^{(2)} + \sum_{j_4=1}^{n_1} W_{i_4 j_4}^{(2)} \sigma \left(\hat{z}_{j_4}^{(1)} \right) \right) \right]$$

Wick = $(\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3})$

$$\times \left\{ C_b^2 + 2C_b C_W \frac{1}{n_1} \sum_{j=1}^{n_1} \mathbb{E} \left[\sigma \left(\hat{z}_j^{(1)} \right) \sigma \left(\hat{z}_j^{(1)} \right) \right] + C_W^2 \frac{1}{n_1^2} \sum_{j_1, j_2=1}^{n_1} \mathbb{E} \left[\sigma \left(\hat{z}_{j_1}^{(1)} \right) \sigma \left(\hat{z}_{j_1}^{(1)} \right) \sigma \left(\hat{z}_{j_2}^{(1)} \right) \sigma \left(\hat{z}_{j_2}^{(1)} \right) \right] \right\}$$

$$= (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3})$$

$$\times \left\{ C_b^2 + 2C_b C_W \langle \sigma(z) \sigma(z) \rangle_{G^{(1)}} + C_W^2 \left[\frac{n_1^2 - 1}{n_1^2} \langle \sigma(z) \sigma(z) \rangle_{G^{(1)}}^2 + \frac{1}{n_1^2} \langle \sigma(z) \sigma(z) \sigma(z) \sigma(z) \rangle_{G^{(1)}} \right] \right\}$$

$$= (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3}) \left\{ \left(\underline{G^{(2)}} \right)^2 + \frac{1}{n_1} C_W^2 \left[\langle \sigma(z) \sigma(z) \sigma(z) \sigma(z) \rangle_{G^{(1)}} - \langle \sigma(z) \sigma(z) \rangle_{G^{(1)}}^2 \right] \right\}$$

Statistics of $\hat{z}_i^{(2)} = b_i^{(2)} + \sum_{j=1}^{n_1} W_{ij}^{(2)} \sigma \left(\hat{z}_j^{(1)} \right)$

$$\mathbb{E} \left[\hat{z}_{i_1}^{(2)} \hat{z}_{i_2}^{(2)} \hat{z}_{i_3}^{(2)} \hat{z}_{i_4}^{(2)} \right]$$

$$= \mathbb{E} \left[\left(b_{i_1}^{(2)} + \sum_{j_1=1}^{n_1} W_{i_1 j_1}^{(2)} \sigma \left(\hat{z}_{j_1}^{(1)} \right) \right) \cdots \left(b_{i_4}^{(2)} + \sum_{j_4=1}^{n_1} W_{i_4 j_4}^{(2)} \sigma \left(\hat{z}_{j_4}^{(1)} \right) \right) \right]$$

Wick = $(\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3})$

$$\times \left\{ C_b^2 + 2C_b C_W \frac{1}{n_1} \sum_{j=1}^{n_1} \mathbb{E} \left[\sigma \left(\hat{z}_j^{(1)} \right) \sigma \left(\hat{z}_j^{(1)} \right) \right] + C_W^2 \frac{1}{n_1^2} \sum_{j_1, j_2=1}^{n_1} \mathbb{E} \left[\sigma \left(\hat{z}_{j_1}^{(1)} \right) \sigma \left(\hat{z}_{j_1}^{(1)} \right) \sigma \left(\hat{z}_{j_2}^{(1)} \right) \sigma \left(\hat{z}_{j_2}^{(1)} \right) \right] \right\}$$

$$= (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3})$$

$$\times \left\{ C_b^2 + 2C_b C_W \langle \sigma(z) \sigma(z) \rangle_{G^{(1)}} + C_W^2 \left[\frac{n_1^2}{n_1^2} \langle \sigma(z) \sigma(z) \rangle_{G^{(1)}} \langle \sigma(z) \sigma(z) \rangle_{G^{(1)}} + \frac{n_1}{n_1^2} \langle \sigma(z) \sigma(z) \sigma(z) \sigma(z) \rangle_{G^{(1)}} \right] \right\}$$

$$= (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3}) \left\{ \left(G^{(2)} \right)^2 + \frac{1}{n_1} C_W^2 \left[\langle \sigma(z) \sigma(z) \sigma(z) \sigma(z) \rangle_{G^{(1)}} - \langle \sigma(z) \sigma(z) \rangle_{G^{(1)}}^2 \right] \right\}$$

Statistics of $\hat{z}_i^{(2)} = b_i^{(2)} + \sum_{j=1}^{n_1} W_{ij}^{(2)} \sigma \left(\hat{z}_j^{(1)} \right)$

$$\begin{aligned}
 & \mathbb{E}[\hat{z}_{i_1}^{(2)} \hat{z}_{i_2}^{(2)} \hat{z}_{i_3}^{(2)} \hat{z}_{i_4}^{(2)}] \Big|_{\text{connected}} \\
 & \equiv \mathbb{E}[\hat{z}_{i_1}^{(2)} \hat{z}_{i_2}^{(2)} \hat{z}_{i_3}^{(2)} \hat{z}_{i_4}^{(2)}] \\
 & \quad - \mathbb{E}[\hat{z}_{i_1}^{(2)} \hat{z}_{i_2}^{(2)}] \mathbb{E}[\hat{z}_{i_3}^{(2)} \hat{z}_{i_4}^{(2)}] - \mathbb{E}[\hat{z}_{i_1}^{(2)} \hat{z}_{i_3}^{(2)}] \mathbb{E}[\hat{z}_{i_2}^{(2)} \hat{z}_{i_4}^{(2)}] - \mathbb{E}[\hat{z}_{i_1}^{(2)} \hat{z}_{i_4}^{(2)}] \mathbb{E}[\hat{z}_{i_2}^{(2)} \hat{z}_{i_3}^{(2)}] \\
 & = \frac{V^{(2)}}{n_1} (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3}) \\
 & \quad \text{with } V^{(2)} = C_W^2 \left[\langle \sigma(z) \sigma(z) \sigma(z) \sigma(z) \rangle_{G^{(1)}} - \langle \sigma(z) \sigma(z) \rangle_{G^{(1)}}^2 \right]
 \end{aligned}$$

Statistics of $\hat{z}_i^{(2)} = b_i^{(2)} + \sum_{j=1}^{n_1} W_{ij}^{(2)} \sigma(\hat{z}_j^{(1)})$

$$\begin{aligned}
& \mathbb{E}[\hat{z}_{i_1}^{(2)} \hat{z}_{i_2}^{(2)} \hat{z}_{i_3}^{(2)} \hat{z}_{i_4}^{(2)}] \Big|_{\text{connected}} \\
& \equiv \mathbb{E}[\hat{z}_{i_1}^{(2)} \hat{z}_{i_2}^{(2)} \hat{z}_{i_3}^{(2)} \hat{z}_{i_4}^{(2)}] \\
& \quad - \mathbb{E}[\hat{z}_{i_1}^{(2)} \hat{z}_{i_2}^{(2)}] \mathbb{E}[\hat{z}_{i_3}^{(2)} \hat{z}_{i_4}^{(2)}] - \mathbb{E}[\hat{z}_{i_1}^{(2)} \hat{z}_{i_3}^{(2)}] \mathbb{E}[\hat{z}_{i_2}^{(2)} \hat{z}_{i_4}^{(2)}] - \mathbb{E}[\hat{z}_{i_1}^{(2)} \hat{z}_{i_4}^{(2)}] \mathbb{E}[\hat{z}_{i_2}^{(2)} \hat{z}_{i_3}^{(2)}] \\
& = \frac{V^{(2)}}{n_1} (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3})
\end{aligned}$$

$$\text{with } V^{(2)} = C_W^2 \left[\langle \sigma(z) \sigma(z) \sigma(z) \sigma(z) \rangle_{G^{(1)}} - \langle \sigma(z) \sigma(z) \rangle_{G^{(1)}}^2 \right]$$

Nearly-Gaussian distribution for $n_1 \gg 1$

[Cf. Gaussian distribution in the first layer:

$$\begin{aligned}
& \mathbb{E}[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)} \hat{z}_{i_3}^{(1)} \hat{z}_{i_4}^{(1)}] \Big|_{\text{connected}} \\
& \equiv \mathbb{E}[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)} \hat{z}_{i_3}^{(1)} \hat{z}_{i_4}^{(1)}] \\
& \quad - \mathbb{E}[\hat{z}_{i_1}^{(1)} \hat{z}_{i_2}^{(1)}] \mathbb{E}[\hat{z}_{i_3}^{(1)} \hat{z}_{i_4}^{(1)}] - \mathbb{E}[\hat{z}_{i_1}^{(1)} \hat{z}_{i_3}^{(1)}] \mathbb{E}[\hat{z}_{i_2}^{(1)} \hat{z}_{i_4}^{(1)}] - \mathbb{E}[\hat{z}_{i_1}^{(1)} \hat{z}_{i_4}^{(1)}] \mathbb{E}[\hat{z}_{i_2}^{(1)} \hat{z}_{i_3}^{(1)}] \\
& = 0
\end{aligned}$$

]

Statistics of $\hat{z}_i^{(2)} = b_i^{(2)} + \sum_{j=1}^{n_1} W_{ij}^{(2)} \sigma \left(\hat{z}_j^{(1)} \right)$

$$\begin{aligned} \mathbb{E} \left[\hat{z}_{i_1}^{(2)} \hat{z}_{i_2}^{(2)} \right] &= G^{(2)} \delta_{i_1 i_2} \\ \mathbb{E} \left[\hat{z}_{i_1}^{(2)} \hat{z}_{i_2}^{(2)} \hat{z}_{i_3}^{(2)} \hat{z}_{i_4}^{(2)} \right] \Big|_{\text{connected}} &= \frac{1}{n_1} V^{(2)} (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3}) \\ \mathbb{E} \left[\hat{z}_{i_1}^{(2)} \hat{z}_{i_2}^{(2)} \hat{z}_{i_3}^{(2)} \hat{z}_{i_4}^{(2)} \hat{z}_{i_5}^{(2)} \hat{z}_{i_6}^{(2)} \right] \Big|_{\text{connected}} &= O \left(\frac{1}{n_1^2} \right) \end{aligned}$$

- Gaussian in the infinite-width limit, too simple; specified by one number (one matrix – kernel – more generally)
- Sparse description at $O(1/n)$; specified by two numbers (two tensors more generally, one of them having four sample indices)
- Interacting neurons at finite width.

Statistics of $\hat{H}_{i_1 i_2}^{(2)}$

$$\hat{H}_{i_1 i_2}^{(2)} \equiv \sum_{\mu, \nu} \lambda_{\mu \nu} \frac{d\hat{z}_{i_1}^{(2)}}{d\theta_{\mu}} \frac{d\hat{z}_{i_2}^{(2)}}{d\theta_{\nu}}$$

Statistics of $\hat{H}_{i_1 i_2}^{(2)}$

$$\begin{aligned}\hat{H}_{i_1 i_2}^{(2)} &\equiv \sum_{\mu, \nu} \lambda_{\mu \nu} \frac{d\hat{z}_{i_1}^{(2)}}{d\theta_{\mu}} \frac{d\hat{z}_{i_2}^{(2)}}{d\theta_{\nu}} \\ &= \lambda_b \sum_{j=1}^{n_2} \frac{d\hat{z}_{i_1}^{(2)}}{db_j^{(2)}} \frac{d\hat{z}_{i_2}^{(2)}}{db_j^{(2)}} + \frac{\lambda_W}{n_1} \sum_{j=1}^{n_2} \sum_{k=1}^{n_1} \frac{d\hat{z}_{i_1}^{(2)}}{dW_{jk}^{(2)}} \frac{d\hat{z}_{i_2}^{(2)}}{dW_{jk}^{(2)}} \\ &\quad + \lambda_b \sum_{j=1}^{n_1} \frac{d\hat{z}_{i_1}^{(2)}}{db_j^{(1)}} \frac{d\hat{z}_{i_2}^{(2)}}{db_j^{(1)}} + \frac{\lambda_W}{n_0} \sum_{j=1}^{n_1} \sum_{k=1}^{n_0} \frac{d\hat{z}_{i_1}^{(2)}}{dW_{jk}^{(1)}} \frac{d\hat{z}_{i_2}^{(2)}}{dW_{jk}^{(1)}}\end{aligned}$$

$$\lambda_{b_{i_1}^{(\ell)} b_{i_2}^{(\ell)}} = \delta_{i_1 i_2} \lambda_b, \quad \lambda_{W_{i_1 j_1}^{(\ell)} W_{i_2 j_2}^{(\ell)}} = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{\lambda_W}{n_{\ell-1}}$$

Statistics of $\hat{H}_{i_1 i_2}^{(2)}$

$$\hat{H}_{i_1 i_2}^{(2)} \equiv \sum_{\mu, \nu} \lambda_{\mu \nu} \frac{d\hat{z}_{i_1}^{(2)}}{d\theta_{\mu}} \frac{d\hat{z}_{i_2}^{(2)}}{d\theta_{\nu}}$$

$$= \lambda_b \sum_{j=1}^{n_2} \frac{d\hat{z}_{i_1}^{(2)}}{db_j^{(2)}} \frac{d\hat{z}_{i_2}^{(2)}}{db_j^{(2)}} + \frac{\lambda_W}{n_1} \sum_{j=1}^{n_2} \sum_{k=1}^{n_1} \frac{d\hat{z}_{i_1}^{(2)}}{dW_{jk}^{(2)}} \frac{d\hat{z}_{i_2}^{(2)}}{dW_{jk}^{(2)}}$$

1st piece

$$+ \lambda_b \sum_{j=1}^{n_1} \frac{d\hat{z}_{i_1}^{(2)}}{db_j^{(1)}} \frac{d\hat{z}_{i_2}^{(2)}}{db_j^{(1)}} + \frac{\lambda_W}{n_0} \sum_{j=1}^{n_1} \sum_{k=1}^{n_0} \frac{d\hat{z}_{i_1}^{(2)}}{dW_{jk}^{(1)}} \frac{d\hat{z}_{i_2}^{(2)}}{dW_{jk}^{(1)}}$$

2nd piece

Statistics of $\hat{H}_{i_1 i_2}^{(2)}$

1st piece, the same as before:

$$\begin{aligned} & \lambda_b \sum_{j=1}^{n_2} \frac{d\hat{z}_{i_1}^{(2)}}{db_j^{(2)}} \frac{d\hat{z}_{i_2}^{(2)}}{db_j^{(2)}} + \frac{\lambda_W}{n_1} \sum_{j=1}^{n_2} \sum_{k=1}^{n_1} \frac{d\hat{z}_{i_1}^{(2)}}{dW_{jk}^{(2)}} \frac{d\hat{z}_{i_2}^{(2)}}{dW_{jk}^{(2)}} \\ &= \lambda_b \sum_{j=1}^{n_2} \delta_{i_1 j} \delta_{i_2 j} + \frac{\lambda_W}{n_1} \sum_{j=1}^{n_2} \sum_{k=1}^{n_1} \delta_{i_1 j} \sigma(\hat{z}_k^{(1)}) \delta_{i_2 j} \sigma(\hat{z}_k^{(1)}) \\ &= \delta_{i_1 i_2} \left\{ \lambda_b + \lambda_W \left[\frac{1}{n_1} \sum_{j=1}^{n_1} \left(\sigma(\hat{z}_j^{(1)}) \right)^2 \right] \right\} \end{aligned}$$

$$\hat{z}_i^{(2)} = b_i^{(2)} + \sum_{j=1}^{n_1} W_{ij}^{(2)} \sigma(\hat{z}_j^{(1)})$$

Statistics of $\hat{H}_{i_1 i_2}^{(2)}$

1st piece, the same as before:

$$\begin{aligned}
 & \lambda_b \sum_{j=1}^{n_2} \frac{d\hat{z}_{i_1}^{(2)}}{db_j^{(2)}} \frac{d\hat{z}_{i_2}^{(2)}}{db_j^{(2)}} + \frac{\lambda_W}{n_1} \sum_{j=1}^{n_2} \sum_{k=1}^{n_1} \frac{d\hat{z}_{i_1}^{(2)}}{dW_{jk}^{(2)}} \frac{d\hat{z}_{i_2}^{(2)}}{dW_{jk}^{(2)}} \\
 &= \lambda_b \sum_{j=1}^{n_2} \delta_{i_1 j} \delta_{i_2 j} + \frac{\lambda_W}{n_1} \sum_{j=1}^{n_2} \sum_{k=1}^{n_1} \delta_{i_1 j} \sigma(\hat{z}_k^{(1)}) \delta_{i_2 j} \sigma(\hat{z}_k^{(1)}) \\
 &= \delta_{i_1 i_2} \left\{ \lambda_b + \lambda_W \left[\frac{1}{n_1} \sum_{j=1}^{n_1} \left(\sigma(\hat{z}_j^{(1)}) \right)^2 \right] \right\}
 \end{aligned}$$

$$\lambda_{W_{i_1 j_1}^{(2)} W_{i_2 j_2}^{(2)}} = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{\lambda_W}{n_1} \quad \text{width-scaling was important.}$$

Statistics of $\hat{H}_{i_1 i_2}^{(2)}$

2nd piece, chain rule:

$$\lambda_b \sum_{j=1}^{n_1} \frac{d\hat{z}_{i_1}^{(2)}}{db_j^{(1)}} \frac{d\hat{z}_{i_2}^{(2)}}{db_j^{(1)}} + \frac{\lambda_W}{n_0} \sum_{j=1}^{n_1} \sum_{k=1}^{n_0} \frac{d\hat{z}_{i_1}^{(2)}}{dW_{jk}^{(1)}} \frac{d\hat{z}_{i_2}^{(2)}}{dW_{jk}^{(1)}}$$

$$= \sum_{m_1, m_2=1}^{n_1} \frac{d\hat{z}_{i_1}^{(2)}}{d\hat{z}_{m_1}^{(1)}} \frac{d\hat{z}_{i_2}^{(2)}}{d\hat{z}_{m_2}^{(1)}} \left[\lambda_b \sum_{j=1}^{n_1} \frac{d\hat{z}_{m_1}^{(1)}}{db_j^{(1)}} \frac{d\hat{z}_{m_2}^{(1)}}{db_j^{(1)}} + \frac{\lambda_W}{n_0} \sum_{j=1}^{n_1} \sum_{k=1}^{n_0} \frac{d\hat{z}_{m_1}^{(2)}}{dW_{jk}^{(1)}} \frac{d\hat{z}_{m_2}^{(2)}}{dW_{jk}^{(1)}} \right]$$

$$\frac{d\hat{z}_i^{(2)}}{d\theta_\mu^{(1)}} = \sum_{m=1}^{n_1} \frac{d\hat{z}_i^{(2)}}{d\hat{z}_m^{(1)}} \frac{d\hat{z}_m^{(1)}}{d\theta_\mu^{(1)}}$$

Statistics of $\hat{H}_{i_1 i_2}^{(2)}$

2nd piece, chain rule:

$$\begin{aligned}
 & \lambda_b \sum_{j=1}^{n_1} \frac{d\hat{z}_{i_1}^{(2)}}{db_j^{(1)}} \frac{d\hat{z}_{i_2}^{(2)}}{db_j^{(1)}} + \frac{\lambda_W}{n_0} \sum_{j=1}^{n_1} \sum_{k=1}^{n_0} \frac{d\hat{z}_{i_1}^{(2)}}{dW_{jk}^{(1)}} \frac{d\hat{z}_{i_2}^{(2)}}{dW_{jk}^{(1)}} \\
 &= \sum_{m_1, m_2=1}^{n_1} \frac{d\hat{z}_{i_1}^{(2)}}{d\hat{z}_{m_1}^{(1)}} \frac{d\hat{z}_{i_2}^{(2)}}{d\hat{z}_{m_2}^{(1)}} \left[\lambda_b \sum_{j=1}^{n_1} \frac{d\hat{z}_{m_1}^{(1)}}{db_j^{(1)}} \frac{d\hat{z}_{m_2}^{(1)}}{db_j^{(1)}} + \frac{\lambda_W}{n_0} \sum_{j=1}^{n_1} \sum_{k=1}^{n_0} \frac{d\hat{z}_{m_1}^{(2)}}{dW_{jk}^{(1)}} \frac{d\hat{z}_{m_2}^{(2)}}{dW_{jk}^{(1)}} \right] \\
 &= \sum_{m_1, m_2=1}^{n_1} \frac{d\hat{z}_{i_1}^{(2)}}{d\hat{z}_{m_1}^{(1)}} \frac{d\hat{z}_{i_2}^{(2)}}{d\hat{z}_{m_2}^{(1)}} \delta_{m_1 m_2} H^{(1)}
 \end{aligned}$$

Statistics of $\hat{H}_{i_1 i_2}^{(2)}$

2nd piece, chain rule:

$$\begin{aligned}
 & \lambda_b \sum_{j=1}^{n_1} \frac{d\hat{z}_{i_1}^{(2)}}{db_j^{(1)}} \frac{d\hat{z}_{i_2}^{(2)}}{db_j^{(1)}} + \frac{\lambda_W}{n_0} \sum_{j=1}^{n_1} \sum_{k=1}^{n_0} \frac{d\hat{z}_{i_1}^{(2)}}{dW_{jk}^{(1)}} \frac{d\hat{z}_{i_2}^{(2)}}{dW_{jk}^{(1)}} \\
 &= \sum_{m_1, m_2=1}^{n_1} \frac{d\hat{z}_{i_1}^{(2)}}{d\hat{z}_{m_1}^{(1)}} \frac{d\hat{z}_{i_2}^{(2)}}{d\hat{z}_{m_2}^{(1)}} \left[\lambda_b \sum_{j=1}^{n_1} \frac{d\hat{z}_{m_1}^{(1)}}{db_j^{(1)}} \frac{d\hat{z}_{m_2}^{(1)}}{db_j^{(1)}} + \frac{\lambda_W}{n_0} \sum_{j=1}^{n_1} \sum_{k=1}^{n_0} \frac{d\hat{z}_{m_1}^{(1)}}{dW_{jk}^{(1)}} \frac{d\hat{z}_{m_2}^{(1)}}{dW_{jk}^{(1)}} \right] \\
 &= \sum_{m_1, m_2=1}^{n_1} \frac{d\hat{z}_{i_1}^{(2)}}{d\hat{z}_{m_1}^{(1)}} \frac{d\hat{z}_{i_2}^{(2)}}{d\hat{z}_{m_2}^{(1)}} \delta_{m_1 m_2} H^{(1)} \\
 &= \sum_{m_1, m_2=1}^{n_1} W_{i_1 m_1}^{(2)} \sigma' \left(\hat{z}_{m_1}^{(1)} \right) W_{i_2 m_2}^{(2)} \sigma' \left(\hat{z}_{m_2}^{(1)} \right) \delta_{m_1 m_2} H^{(1)}
 \end{aligned}$$

$$\hat{z}_i^{(2)} = b_i^{(2)} + \sum_{j=1}^{n_1} W_{ij}^{(2)} \sigma \left(\hat{z}_j^{(1)} \right)$$

Statistics of $\hat{H}_{i_1 i_2}^{(2)}$

2nd piece, chain rule:

$$\begin{aligned}
 & \lambda_b \sum_{j=1}^{n_1} \frac{d\hat{z}_{i_1}^{(2)}}{db_j^{(1)}} \frac{d\hat{z}_{i_2}^{(2)}}{db_j^{(1)}} + \frac{\lambda_W}{n_0} \sum_{j=1}^{n_1} \sum_{k=1}^{n_0} \frac{d\hat{z}_{i_1}^{(2)}}{dW_{jk}^{(1)}} \frac{d\hat{z}_{i_2}^{(2)}}{dW_{jk}^{(1)}} \\
 &= \sum_{m_1, m_2=1}^{n_1} \frac{d\hat{z}_{i_1}^{(2)}}{d\hat{z}_{m_1}^{(1)}} \frac{d\hat{z}_{i_2}^{(2)}}{d\hat{z}_{m_2}^{(1)}} \left[\lambda_b \sum_{j=1}^{n_1} \frac{d\hat{z}_{m_1}^{(1)}}{db_j^{(1)}} \frac{d\hat{z}_{m_2}^{(1)}}{db_j^{(1)}} + \frac{\lambda_W}{n_0} \sum_{j=1}^{n_1} \sum_{k=1}^{n_0} \frac{d\hat{z}_{m_1}^{(1)}}{dW_{jk}^{(1)}} \frac{d\hat{z}_{m_2}^{(1)}}{dW_{jk}^{(1)}} \right] \\
 &= \sum_{m_1, m_2=1}^{n_1} \frac{d\hat{z}_{i_1}^{(2)}}{d\hat{z}_{m_1}^{(1)}} \frac{d\hat{z}_{i_2}^{(2)}}{d\hat{z}_{m_2}^{(1)}} \delta_{m_1 m_2} H^{(1)} \\
 &= \sum_{m_1, m_2=1}^{n_1} W_{i_1 m_1}^{(2)} \sigma' \left(\hat{z}_{m_1}^{(1)} \right) W_{i_2 m_2}^{(2)} \sigma' \left(\hat{z}_{m_2}^{(1)} \right) \delta_{m_1 m_2} H^{(1)} \\
 &= \sum_{m=1}^{n_1} W_{i_1 m}^{(2)} W_{i_2 m}^{(2)} \sigma' \left(\hat{z}_m^{(1)} \right) \sigma' \left(\hat{z}_m^{(1)} \right) H^{(1)}
 \end{aligned}$$

Statistics of $\hat{H}_{i_1 i_2}^{(2)}$

Putting things together, NTK *forward* equation:

$$\hat{H}_{i_1 i_2}^{(2)} = \delta_{i_1 i_2} \left\{ \lambda_b + \lambda_W \left[\frac{1}{n_1} \sum_{j=1}^{n_1} \left(\sigma(\hat{z}_j^{(1)}) \right)^2 \right] \right\} + \sum_{m=1}^{n_1} W_{i_1 m}^{(2)} W_{i_2 m}^{(2)} \sigma'(\hat{z}_m^{(1)}) \sigma'(\hat{z}_m^{(1)}) H^{(1)}$$

Statistics of $\hat{H}_{i_1 i_2}^{(2)}$

Putting things together, NTK *forward* equation:

$$\hat{H}_{i_1 i_2}^{(2)} = \delta_{i_1 i_2} \left\{ \lambda_b + \lambda_W \left[\frac{1}{n_1} \sum_{j=1}^{n_1} \left(\sigma(\hat{z}_j^{(1)}) \right)^2 \right] \right\} + \sum_{m=1}^{n_1} W_{i_1 m}^{(2)} W_{i_2 m}^{(2)} \sigma'(\hat{z}_m^{(1)}) \sigma'(\hat{z}_m^{(1)}) H^{(1)}$$

- “Stochastic”: it fluctuates from instantiation to instantiation.
- “Defrosted”: it *can* evolve during training.

Statistics of $\hat{H}_{i_1 i_2}^{(2)}$

Putting things together, NTK *forward* equation:

$$\hat{H}_{i_1 i_2}^{(2)} = \delta_{i_1 i_2} \left\{ \lambda_b + \lambda_W \left[\frac{1}{n_1} \sum_{j=1}^{n_1} \left(\sigma(\hat{z}_j^{(1)}) \right)^2 \right] \right\} + \sum_{m=1}^{n_1} W_{i_1 m}^{(2)} W_{i_2 m}^{(2)} \sigma'(\hat{z}_m^{(1)}) \sigma'(\hat{z}_m^{(1)}) H^{(1)}$$

Fun for the weekend (solutions in §8):

$$\mathbb{E} \left[\hat{H}_{i_1 i_2}^{(2)} \right] = \delta_{i_1 i_2} \left[\lambda_b + \lambda_W \langle \sigma(z) \sigma(z) \rangle_{G^{(1)}} + C_W H^{(1)} \langle \sigma'(z) \sigma'(z) \rangle_{G^{(1)}} \right] \equiv \delta_{i_1 i_2} H^{(2)}$$

$$\mathbb{E} \left[\hat{H}_{i_1 i_2}^{(2)} \hat{H}_{i_3 i_4}^{(2)} \right] - \mathbb{E} \left[\hat{H}_{i_1 i_2}^{(2)} \right] \mathbb{E} \left[\hat{H}_{i_3 i_4}^{(2)} \right] = \frac{1}{n_1} \left[\delta_{i_1 i_2} \delta_{i_3 i_4} A^{(2)} + (\delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3}) B^{(2)} \right] = O(1/n)$$

$$\mathbb{E} \left[\hat{H}_{i_1 i_2}^{(2)} \hat{z}_{i_3}^{(2)} \hat{z}_{i_4}^{(2)} \right] - \mathbb{E} \left[\hat{H}_{i_1 i_2}^{(2)} \right] \mathbb{E} \left[\hat{z}_{i_3}^{(2)} \hat{z}_{i_4}^{(2)} \right] = \frac{1}{n_1} \left[\delta_{i_1 i_2} \delta_{i_3 i_4} D^{(2)} + (\delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3}) F^{(2)} \right] = O(1/n)$$

Statistics of $\hat{H}_{i_1 i_2}^{(2)}$ and beyond

Putting things together, NTK *forward* equation:

$$\hat{H}_{i_1 i_2}^{(2)} = \delta_{i_1 i_2} \left\{ \lambda_b + \lambda_W \left[\frac{1}{n_1} \sum_{j=1}^{n_1} \left(\sigma(\hat{z}_j^{(1)}) \right)^2 \right] \right\} + \sum_{m=1}^{n_1} W_{i_1 m}^{(2)} W_{i_2 m}^{(2)} \sigma'(\hat{z}_m^{(1)}) \sigma'(\hat{z}_m^{(1)}) H^{(1)}$$

Fun for the weekend (solutions in §8, §11.2, and §∞.3):

$$\mathbb{E} \left[\hat{H}_{i_1 i_2}^{(2)} \right] = \delta_{i_1 i_2} \left[\lambda_b + \lambda_W \langle \sigma(z) \sigma(z) \rangle_{G^{(1)}} + C_W H^{(1)} \langle \sigma'(z) \sigma'(z) \rangle_{G^{(1)}} \right] \equiv \delta_{i_1 i_2} H^{(2)}$$

$$\mathbb{E} \left[\hat{H}_{i_1 i_2}^{(2)} \hat{H}_{i_3 i_4}^{(2)} \right] - \mathbb{E} \left[\hat{H}_{i_1 i_2}^{(2)} \right] \mathbb{E} \left[\hat{H}_{i_3 i_4}^{(2)} \right] = \frac{1}{n_1} \left[\delta_{i_1 i_2} \delta_{i_3 i_4} A^{(2)} + (\delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3}) B^{(2)} \right] = O(1/n)$$

$$\mathbb{E} \left[\hat{H}_{i_1 i_2}^{(2)} \hat{z}_{i_3}^{(2)} \hat{z}_{i_4}^{(2)} \right] - \mathbb{E} \left[\hat{H}_{i_1 i_2}^{(2)} \right] \mathbb{E} \left[\hat{z}_{i_3}^{(2)} \hat{z}_{i_4}^{(2)} \right] = \frac{1}{n_1} \left[\delta_{i_1 i_2} \delta_{i_3 i_4} D^{(2)} + (\delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3}) F^{(2)} \right] = O(1/n)$$

$$\mathbb{E} \left[\widehat{\mathrm{d}H}_{i_0 i_1 i_2}^{(2)} \hat{z}_{i_3}^{(2)} \right] = \frac{1}{n_1} \left[\delta_{i_0 i_3} \delta_{i_1 i_2} P^{(2)} + (\delta_{i_0 i_1} \delta_{i_2 i_3} + \delta_{i_0 i_2} \delta_{i_1 i_3}) Q^{(2)} \right] = O(1/n)$$

$$\mathbb{E} \left[\widehat{\mathrm{d}dH}^{(2)} \right] = O(1/n) \quad \text{[*for smooth activation functions]}$$

Statistics of $\hat{H}_{i_1 i_2}^{(2)}$ and beyond

Putting things together, NTK *forward* equation:

$$\hat{H}_{i_1 i_2}^{(2)} = \delta_{i_1 i_2} \left\{ \lambda_b + \lambda_W \left[\frac{1}{n_1} \sum_{j=1}^{n_1} \left(\sigma(\hat{z}_j^{(1)}) \right)^2 \right] \right\} + \sum_{m=1}^{n_1} W_{i_1 m}^{(2)} W_{i_2 m}^{(2)} \sigma'(\hat{z}_m^{(1)}) \sigma'(\hat{z}_m^{(1)}) H^{(1)}$$

Fun for the weekend (solutions in §8, §11.2, and §∞.3):

$$\mathbb{E} \left[\hat{H}_{i_1 i_2}^{(2)} \right] = \delta_{i_1 i_2} \left[\lambda_b + \lambda_W \langle \sigma(z) \sigma(z) \rangle_{G^{(1)}} + C_W H^{(1)} \langle \sigma'(z) \sigma'(z) \rangle_{G^{(1)}} \right] \equiv \delta_{i_1 i_2} H^{(2)}$$

$$\mathbb{E} \left[\hat{H}_{i_1 i_2}^{(2)} \hat{H}_{i_3 i_4}^{(2)} \right] - \mathbb{E} \left[\hat{H}_{i_1 i_2}^{(2)} \right] \mathbb{E} \left[\hat{H}_{i_3 i_4}^{(2)} \right] = \frac{1}{n_1} \left[\delta_{i_1 i_2} \delta_{i_3 i_4} A^{(2)} + (\delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3}) B^{(2)} \right] = O(1/n)$$

$$\mathbb{E} \left[\hat{H}_{i_1 i_2}^{(2)} \hat{z}_{i_3}^{(2)} \hat{z}_{i_4}^{(2)} \right] - \mathbb{E} \left[\hat{H}_{i_1 i_2}^{(2)} \right] \mathbb{E} \left[\hat{z}_{i_3}^{(2)} \hat{z}_{i_4}^{(2)} \right] = \frac{1}{n_1} \left[\delta_{i_1 i_2} \delta_{i_3 i_4} D^{(2)} + (\delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3}) F^{(2)} \right] = O(1/n)$$

$$\mathbb{E} \left[\widehat{\mathrm{d}H}_{i_0 i_1 i_2}^{(2)} \hat{z}_{i_3}^{(2)} \right] = \frac{1}{n_1} \left[\delta_{i_0 i_3} \delta_{i_1 i_2} P^{(2)} + (\delta_{i_0 i_1} \delta_{i_2 i_3} + \delta_{i_0 i_2} \delta_{i_1 i_3}) Q^{(2)} \right] = O(1/n)$$

$$\mathbb{E} \left[\widehat{\mathrm{d}dH}^{(2)} \right] = O(1/n) \quad \text{[*for smooth activation functions]}$$

Representation Learning

Statistics of $\widehat{H}_{i_1 i_2}^{(2)}$ and beyond

$$z_{i;\delta}^* = \widehat{z}_{i;\delta} - \sum \widehat{H}_{ij;\delta\tilde{\alpha}_1} \left(\widehat{H}^{-1} \right)^{jk;\tilde{\alpha}_1\tilde{\alpha}_2} [\widehat{z}_{k;\tilde{\alpha}} - y_{k;\tilde{\alpha}}] \\ + \text{despicable}(y, \widehat{z}, \widehat{H}, \widehat{\mathrm{d}H}, \widehat{\mathrm{d}dH}; \text{algorithm})$$

$$H^* \neq \widehat{H}$$

$$\mathbb{E} \left[\widehat{\mathrm{d}H}_{i_0 i_1 i_2}^{(2)} \widehat{z}_{i_3}^{(2)} \right] = \frac{1}{n_1} \left[\delta_{i_0 i_3} \delta_{i_1 i_2} P^{(2)} + (\delta_{i_0 i_1} \delta_{i_2 i_3} + \delta_{i_0 i_2} \delta_{i_1 i_3}) Q^{(2)} \right] = O(1/n) \\ \mathbb{E} \left[\widehat{\mathrm{d}dH}^{(2)} \right] = O(1/n) \quad [* \text{for smooth activation functions}]$$

Representation Learning

Statistics of Two-Layer Neural Networks

- Two interpretations:
 - (i) outputs, NTK, ... of a two-layer network; or
 - (ii) preactivations, mid-layer NTK, ... in the second layer of a deeper network.
- Neurons *do* talk to each other; they are statistically *dependent*.
- Yes representation learning (and yes algorithm dependence); they can now capture rich dynamics of real, finite-width, neural networks.

Statistics of Two-Layer Neural Networks

- Two interpretations:
 - (i) outputs, NTK, ... of a two-layer network; or
 - (ii) preactivations, mid-layer NTK, ... in the second layer of a deeper network.
- Neurons *do* talk to each other; they are statistically *dependent*.
- Yes representation learning (and yes algorithm dependence); they can now capture rich dynamics of real, finite-width, neural networks.

But what is being amplified by deep learning?

4. Deep Neural Networks

$$p(\hat{z}^{(\ell)}, \hat{H}^{(\ell)}, \widehat{dH}^{(\ell)}, \dots)$$

Statistics of $\hat{z}_i^{(\ell+1)} = b_i^{(\ell+1)} + \sum_{j=1}^{n_1} W_{ij}^{(\ell+1)} \sigma \left(\hat{z}_j^{(\ell)} \right)$

$$\mathbb{E} \left[\hat{z}_{i_1}^{(\ell+1)} \hat{z}_{i_2}^{(\ell+1)} \right] = \mathbb{E} \left[\left(b_{i_1}^{(\ell+1)} + \sum_{j_1=1}^{n_\ell} W_{i_1 j_1}^{(\ell+1)} \sigma \left(\hat{z}_{j_1}^{(\ell)} \right) \right) \left(b_{i_2}^{(\ell+1)} + \sum_{j_2=1}^{n_\ell} W_{i_2 j_2}^{(\ell+1)} \sigma \left(\hat{z}_{j_2}^{(\ell)} \right) \right) \right]$$

$$\mathbb{E} \left[b_{i_1}^{(\ell+1)} b_{i_2}^{(\ell+1)} \right] = \delta_{i_1 i_2} C_b, \quad \mathbb{E} \left[W_{i_1 j_1}^{(\ell+1)} W_{i_2 j_2}^{(\ell+1)} \right] = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{C_W}{n_\ell}$$

Statistics of $\hat{z}_i^{(\ell+1)} = b_i^{(\ell+1)} + \sum_{j=1}^{n_1} W_{ij}^{(\ell+1)} \sigma \left(\hat{z}_j^{(\ell)} \right)$

$$\mathbb{E} \left[\hat{z}_{i_1}^{(\ell+1)} \hat{z}_{i_2}^{(\ell+1)} \right] = \mathbb{E} \left[\left(b_{i_1}^{(\ell+1)} + \sum_{j_1=1}^{n_\ell} W_{i_1 j_1}^{(\ell+1)} \sigma \left(\hat{z}_{j_1}^{(\ell)} \right) \right) \left(b_{i_2}^{(\ell+1)} + \sum_{j_2=1}^{n_\ell} W_{i_2 j_2}^{(\ell+1)} \sigma \left(\hat{z}_{j_2}^{(\ell)} \right) \right) \right]$$

$$\text{Wick} = C_b \delta_{i_1 i_2} + \sum_{j_1, j_2=1}^{n_\ell} \frac{C_W}{n_\ell} \delta_{i_1 i_2} \delta_{j_1 j_2} \mathbb{E} \left[\sigma \left(\hat{z}_{j_1}^{(\ell)} \right) \sigma \left(\hat{z}_{j_2}^{(\ell)} \right) \right]$$

$$\begin{aligned} \text{arrange} &= \delta_{i_1 i_2} \left[C_b + C_W \left(\frac{1}{n_\ell} \sum_{j=1}^{n_\ell} \mathbb{E} \left[\sigma \left(\hat{z}_j^{(\ell)} \right) \sigma \left(\hat{z}_j^{(\ell)} \right) \right] \right) \right] \\ &= \delta_{i_1 i_2} \left[C_b + C_W \langle \sigma(z) \sigma(z) \rangle_{G^{(\ell)}} + O \left(\frac{1}{n} \right) \right] \equiv \delta_{i_1 i_2} G^{(\ell+1)} \end{aligned}$$

$$\mathbb{E} \left[b_{i_1}^{(\ell+1)} b_{i_2}^{(\ell+1)} \right] = \delta_{i_1 i_2} C_b, \quad \mathbb{E} \left[W_{i_1 j_1}^{(\ell+1)} W_{i_2 j_2}^{(\ell+1)} \right] = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{C_W}{n_\ell}$$

Statistics of $\hat{z}_i^{(\ell+1)} = b_i^{(\ell+1)} + \sum_{j=1}^{n_1} W_{ij}^{(\ell+1)} \sigma \left(\hat{z}_j^{(\ell)} \right)$

$$\begin{aligned}
 \mathbb{E} \left[\hat{z}_{i_1}^{(\ell+1)} \hat{z}_{i_2}^{(\ell+1)} \right] &= \mathbb{E} \left[\left(b_{i_1}^{(\ell+1)} + \sum_{j_1=1}^{n_\ell} W_{i_1 j_1}^{(\ell+1)} \sigma \left(\hat{z}_{j_1}^{(\ell)} \right) \right) \left(b_{i_2}^{(\ell+1)} + \sum_{j_2=1}^{n_\ell} W_{i_2 j_2}^{(\ell+1)} \sigma \left(\hat{z}_{j_2}^{(\ell)} \right) \right) \right] \\
 &= C_b \delta_{i_1 i_2} + \sum_{j_1, j_2=1}^{n_\ell} \frac{C_W}{n_\ell} \delta_{i_1 i_2} \delta_{j_1 j_2} \mathbb{E} \left[\sigma \left(\hat{z}_{j_1}^{(\ell)} \right) \sigma \left(\hat{z}_{j_2}^{(\ell)} \right) \right] \\
 &= \delta_{i_1 i_2} \left[C_b + C_W \left(\frac{1}{n_\ell} \sum_{j=1}^{n_\ell} \mathbb{E} \left[\sigma \left(\hat{z}_j^{(\ell)} \right) \sigma \left(\hat{z}_j^{(\ell)} \right) \right] \right) \right] \\
 \text{leading} \quad &= \delta_{i_1 i_2} \left[C_b + C_W \langle \sigma(z) \sigma(z) \rangle_{G^{(\ell)}} + O \left(\frac{1}{n} \right) \right] \equiv \delta_{i_1 i_2} G^{(\ell+1)}
 \end{aligned}$$

Statistics of $\hat{z}_i^{(\ell+1)} = b_i^{(\ell+1)} + \sum_{j=1}^{n_1} W_{ij}^{(\ell+1)} \sigma \left(\hat{z}_j^{(\ell)} \right)$

Two-point:

$$G^{(\ell+1)} = C_b + C_W \langle \sigma(z) \sigma(z) \rangle_{G^{(\ell)}} + O\left(\frac{1}{n}\right)$$

Statistics of $\hat{z}_i^{(\ell+1)} = b_i^{(\ell+1)} + \sum_{j=1}^{n_1} W_{ij}^{(\ell+1)} \sigma \left(\hat{z}_j^{(\ell)} \right)$

Two-point:

$$G^{(\ell+1)} = C_b + C_W \langle \sigma(z) \sigma(z) \rangle_{G^{(\ell)}} + O \left(\frac{1}{n} \right)$$

Four-point:

$$\begin{aligned} \frac{1}{n_\ell} V^{(\ell+1)} = & \frac{1}{n_\ell} C_W^2 \left[\langle \sigma(z) \sigma(z) \sigma(z) \sigma(z) \rangle_{G^{(\ell)}} - \langle \sigma(z) \sigma(z) \rangle_{G^{(\ell)}}^2 \right] \\ & + \frac{C_W^2}{4n_{\ell-1}} \frac{V^{(\ell)}}{(G^{(\ell)})^4} \left\langle \sigma(z) \sigma(z) \left(z^2 - G^{(\ell)} \right) \right\rangle_{G^{(\ell)}}^2 + O \left(\frac{1}{n^2} \right) \end{aligned}$$

Statistics of $\hat{H}_{i_1 i_2}^{(\ell+1)}$

Two-point:

$$G^{(\ell+1)} = C_b + C_W \langle \sigma(z) \sigma(z) \rangle_{G^{(\ell)}} + O\left(\frac{1}{n}\right)$$

Four-point:

$$\begin{aligned} \frac{1}{n_\ell} V^{(\ell+1)} = & \frac{1}{n_\ell} C_W^2 \left[\langle \sigma(z) \sigma(z) \sigma(z) \sigma(z) \rangle_{G^{(\ell)}} - \langle \sigma(z) \sigma(z) \rangle_{G^{(\ell)}}^2 \right] \\ & + \frac{C_W^2}{4n_{\ell-1}} \frac{V^{(\ell)}}{(G^{(\ell)})^4} \left\langle \sigma(z) \sigma(z) \left(z^2 - G^{(\ell)} \right) \right\rangle_{G^{(\ell)}}^2 + O\left(\frac{1}{n^2}\right) \end{aligned}$$

NTK mean:

$$H^{(\ell+1)} = \lambda_b + \lambda_W \langle \sigma(z) \sigma(z) \rangle_{G^{(\ell)}} + C_W H^{(\ell)} \langle \sigma'(z) \sigma'(z) \rangle_{G^{(\ell)}} + O\left(\frac{1}{n}\right)$$

Statistics of $\hat{H}_{i_1 i_2}^{(\ell+1)}$

NTK *forward* equation:

$$\begin{aligned}
 \hat{H}_{i_1 i_2}^{(\ell+1)} &\equiv \sum_{\mu, \nu} \lambda_{\mu \nu} \frac{d\hat{z}_{i_1}^{(\ell+1)}}{d\theta_{\mu}} \frac{d\hat{z}_{i_2}^{(\ell+1)}}{d\theta_{\nu}} \\
 &= \lambda_b \sum_{j=1}^{n_{\ell+1}} \frac{d\hat{z}_{i_1}^{(\ell+1)}}{db_j^{(\ell+1)}} \frac{d\hat{z}_{i_2}^{(\ell+1)}}{db_j^{(\ell+1)}} + \frac{\lambda_W}{n_{\ell}} \sum_{j=1}^{n_{\ell+1}} \sum_{k=1}^{n_{\ell}} \frac{d\hat{z}_{i_1}^{(\ell+1)}}{dW_{jk}^{(\ell+1)}} \frac{d\hat{z}_{i_2}^{(\ell+1)}}{dW_{jk}^{(\ell+1)}} && \text{1st trivial piece} \\
 &\quad + \sum_{m_1, m_2=1}^{n_{\ell}} \frac{d\hat{z}_{i_1}^{(\ell+1)}}{d\hat{z}_{m_1}^{(\ell)}} \frac{d\hat{z}_{i_2}^{(\ell+1)}}{d\hat{z}_{m_2}^{(\ell)}} \hat{H}_{m_1 m_2}^{(\ell)} && \text{2nd chain-rule piece}
 \end{aligned}$$

Statistics of $\hat{H}_{i_1 i_2}^{(\ell+1)}$

NTK *forward* equation:

$$\begin{aligned}\hat{H}_{i_1 i_2}^{(\ell+1)} &\equiv \sum_{\mu, \nu} \lambda_{\mu \nu} \frac{d\hat{z}_{i_1}^{(\ell+1)}}{d\theta_{\mu}} \frac{d\hat{z}_{i_2}^{(\ell+1)}}{d\theta_{\nu}} \\ &= \delta_{i_1 i_2} \left\{ \lambda_b + \lambda_W \left[\frac{1}{n_{\ell}} \sum_{j=1}^{n_{\ell}} \left(\sigma(\hat{z}_j^{(\ell)}) \right)^2 \right] \right\} \\ &\quad + \sum_{m_1, m_2=1}^{n_{\ell}} W_{i_1 m_1}^{(\ell+1)} W_{i_2 m_2}^{(\ell+1)} \sigma'(\hat{z}_{m_1}^{(\ell)}) \sigma'(\hat{z}_{m_2}^{(\ell)}) \hat{H}_{m_1 m_2}^{(\ell)}\end{aligned}$$

1st trivial piece

2nd chain-rule piece

Statistics of $\hat{H}_{i_1 i_2}^{(\ell+1)}$ and beyond

Two-point:

$$G^{(\ell+1)} = C_b + C_W \langle \sigma(z) \sigma(z) \rangle_{G^{(\ell)}} + O\left(\frac{1}{n}\right)$$

Four-point:

$$\begin{aligned} \frac{1}{n_\ell} V^{(\ell+1)} = & \frac{1}{n_\ell} C_W^2 \left[\langle \sigma(z) \sigma(z) \sigma(z) \sigma(z) \rangle_{G^{(\ell)}} - \langle \sigma(z) \sigma(z) \rangle_{G^{(\ell)}}^2 \right] \\ & + \frac{C_W^2}{4n_{\ell-1}} \frac{V^{(\ell)}}{(G^{(\ell)})^4} \left\langle \sigma(z) \sigma(z) \left(z^2 - G^{(\ell)} \right) \right\rangle_{G^{(\ell)}}^2 + O\left(\frac{1}{n^2}\right) \end{aligned}$$

NTK mean:

$$H^{(\ell+1)} = \lambda_b + \lambda_W \langle \sigma(z) \sigma(z) \rangle_{G^{(\ell)}} + C_W H^{(\ell)} \langle \sigma'(z) \sigma'(z) \rangle_{G^{(\ell)}} + O\left(\frac{1}{n}\right)$$

NTK fluctuations (§8)

ddNTK (§∞.3)

Similarly for $A^{(\ell)}, B^{(\ell)}, D^{(\ell)}, F^{(\ell)}, \underbrace{P^{(\ell)}, Q^{(\ell)}}_{\text{dNTK (§11.2)}}, \underbrace{R^{(\ell)}, S^{(\ell)}, T^{(\ell)}, U^{(\ell)}}_{\text{ddNTK (§∞.3)}}$

Statistics of $\hat{H}_{i_1 i_2}^{(\ell+1)}$ and beyond

Two-point:

$$G^{(\ell+1)} = C_b + C_W \langle \sigma(z) \sigma(z) \rangle_{G^{(\ell)}} + O\left(\frac{1}{n}\right)$$

Four-point:

$$\begin{aligned} \frac{1}{n_\ell} V^{(\ell+1)} = & \frac{1}{n_\ell} C_W^2 \left[\langle \sigma(z) \sigma(z) \sigma(z) \sigma(z) \rangle_{G^{(\ell)}} - \langle \sigma(z) \sigma(z) \rangle_{G^{(\ell)}}^2 \right] \\ & + \frac{C_W^2}{4n_{\ell-1}} \frac{V^{(\ell)}}{(G^{(\ell)})^4} \left\langle \sigma(z) \sigma(z) \left(z^2 - G^{(\ell)} \right) \right\rangle_{G^{(\ell)}}^2 + O\left(\frac{1}{n^2}\right) \end{aligned} \quad O\left(\frac{1}{n}\right)$$

NTK mean:

$$H^{(\ell+1)} = \lambda_b + \lambda_W \langle \sigma(z) \sigma(z) \rangle_{G^{(\ell)}} + C_W H^{(\ell)} \langle \sigma'(z) \sigma'(z) \rangle_{G^{(\ell)}} + O\left(\frac{1}{n}\right)$$

NTK fluctuations (§8)

ddNTK (§∞.3)

Similarly for $A^{(\ell)}, B^{(\ell)}, D^{(\ell)}, F^{(\ell)}, \underbrace{P^{(\ell)}, Q^{(\ell)}}_{\text{dNTK (§11.2)}}, \underbrace{R^{(\ell)}, S^{(\ell)}, T^{(\ell)}, U^{(\ell)}}_{\text{ddNTK (§∞.3)}}$ $O\left(\frac{1}{n}\right)$

Statistics of $\hat{H}_{i_1 i_2}^{(\ell+1)}$ and beyond

Two-point:

$$G^{(\ell+1)} = C_b + C_W \langle \sigma(z) \sigma(z) \rangle_{G^{(\ell)}} + O\left(\frac{1}{n}\right)$$

Four-point:

$$\begin{aligned} \frac{1}{n_\ell} V^{(\ell+1)} = & \frac{1}{n_\ell} C_W^2 \left[\langle \sigma(z) \sigma(z) \sigma(z) \sigma(z) \rangle_{G^{(\ell)}} - \langle \sigma(z) \sigma(z) \rangle_{G^{(\ell)}}^2 \right] \\ & + \frac{C_W^2}{4n_{\ell-1}} \frac{V^{(\ell)}}{(G^{(\ell)})^4} \left\langle \sigma(z) \sigma(z) \left(z^2 - G^{(\ell)} \right) \right\rangle_{G^{(\ell)}}^2 + O\left(\frac{1}{n^2}\right) \end{aligned} \quad O\left(\frac{L}{n}\right)$$

NTK mean:

$$H^{(\ell+1)} = \lambda_b + \lambda_W \langle \sigma(z) \sigma(z) \rangle_{G^{(\ell)}} + C_W H^{(\ell)} \langle \sigma'(z) \sigma'(z) \rangle_{G^{(\ell)}} + O\left(\frac{1}{n}\right)$$

NTK fluctuations (§8)

ddNTK (§∞.3)

Similarly for $A^{(\ell)}, B^{(\ell)}, D^{(\ell)}, F^{(\ell)}, \underbrace{P^{(\ell)}, Q^{(\ell)}}_{\text{dNTK (§11.2)}}, \underbrace{R^{(\ell)}, S^{(\ell)}, T^{(\ell)}, U^{(\ell)}}_{\text{ddNTK (§∞.3)}}$ $O\left(\frac{L}{n}\right)$

The Principle of Sparsity for WIDE Neural Networks

$$\overset{\text{Sho}}{p(\theta)} \xrightarrow{\quad} p\left(\hat{z}, \hat{H}, \widehat{\mathrm{d}H}, \dots\right) \xrightarrow{\quad} p(z^*)$$

statistics at *initialization* statistics *after training*

The Principle of Sparsity for WIDE Neural Networks

$$p(\theta) \rightarrow p\left(\widehat{z}, \widehat{H}, \widehat{\mathrm{d}H}, \dots\right) \rightarrow p(z^{\star})$$

statistics at *initialization*

statistics *after training*

- Infinite width:

$$p\left(\widehat{z}, \widehat{H}\right) \text{ specified by } G^{(L)}, H^{(L)}$$

The Principle of Sparsity for WIDE Neural Networks

$$p(\theta) \rightarrow p\left(\widehat{z}, \widehat{H}, \widehat{\mathrm{d}H}, \dots\right) \rightarrow p(z^\star)$$

statistics at *initialization*

statistics *after training*

- Infinite width:

$$p\left(\widehat{z}, \widehat{H}\right) \text{ specified by } G^{(L)}, H^{(L)}$$

- Large-but-finite width at $O\left(\frac{1}{n}\right)$ [$n_1, n_2, \dots, n_{L-1} \gg L$]:

$$p\left(\widehat{z}, \widehat{H}, \widehat{\mathrm{d}H}, \widehat{\mathrm{d}\mathrm{d}H}\right) \text{ specified by } G^{(L)}, H^{(L)}, V^{(L)}, A^{(L)}, B^{(L)}, D^{(L)}, F^{(L)}, P^{(L)}, Q^{(L)}, R^{(L)}, S^{(L)}, T^{(L)}, U^{(L)}$$

All determined through recursion relations
(RG-flow interpretation: §4.6)

$$G^{(L)}, H^{(L)}, V^{(L)}, A^{(L)}, B^{(L)}, D^{(L)}, F^{(L)}, P^{(L)}, Q^{(L)}, R^{(L)}, S^{(L)}, T^{(L)}, U^{(L)}$$

All determined through recursion relations
(RG-flow interpretation: §4.6)

Next Lecture: Solving Recursions “The Principle of Criticality” for DEEP Neural Networks

$$G^{(L)}, H^{(L)}, V^{(L)}, A^{(L)}, B^{(L)}, D^{(L)}, F^{(L)}, P^{(L)}, Q^{(L)}, R^{(L)}, S^{(L)}, T^{(L)}, U^{(L)}$$

One more thing...

$$\mathbb{E}[\hat{z}_{i_1}^{(\ell)} \hat{z}_{i_2}^{(\ell)}] = \delta_{i_1 i_2} G^{(\ell)} = \delta_{i_1 i_2} \left[K^{(\ell)} + O\left(\frac{1}{n}\right) \right]$$

$$K^{(\ell+1)} = C_b + C_W \langle \sigma(z) \sigma(z) \rangle_{K^{(\ell)}}$$

$$\mathbb{E}[\hat{H}_{i_1 i_2}^{(\ell)}] = \delta_{i_1 i_2} H^{(\ell)} = \delta_{i_1 i_2} \left[\Theta^{(\ell)} + O\left(\frac{1}{n}\right) \right]$$

$$\Theta^{(\ell+1)} = \lambda_b + \lambda_W \langle \sigma(z) \sigma(z) \rangle_{K^{(\ell)}} + C_W \Theta^{(\ell)} \langle \sigma'(z) \sigma'(z) \rangle_{K^{(\ell)}}$$

$$\mathbb{E} \left[\hat{z}_{i_1}^{(\ell)} \hat{z}_{i_2}^{(\ell)} \hat{z}_{i_3}^{(\ell)} \hat{z}_{i_4}^{(\ell)} \right] \Big|_{\text{connected}} = \frac{1}{n_{\ell-1}} V^{(\ell)} (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3})$$

$$\begin{aligned} \frac{1}{n_\ell} V^{(\ell+1)} = & \frac{1}{n_\ell} C_W^2 \left[\langle \sigma(z) \sigma(z) \sigma(z) \sigma(z) \rangle_{K^{(\ell)}} - \langle \sigma(z) \sigma(z) \rangle_{K^{(\ell)}}^2 \right] \\ & + \frac{C_W^2}{4n_{\ell-1}} \frac{V^{(\ell)}}{(K^{(\ell)})^4} \left\langle \sigma(z) \sigma(z) \left(z^2 - K^{(\ell)} \right) \right\rangle_{K^{(\ell)}}^2 + O\left(\frac{1}{n^2}\right) \end{aligned}$$