

Effective Theory of Deep Learning

Beyond the Infinite-Width Limit

Dan Roberts^a and Sho Yaida^b

^aMIT, IAIFI, & Salesforce, ^bFacebook AI Research

Deep Learning Theory Summer School at Princeton
July 27, 2021 – August 8, 2021

Course Plan

~~Lecture 1~~ **Initialization, Linear Models**

▶ ~~§0 + §7.1 + §10.4~~

~~Lecture 2~~ **Quadratic Models & Nearly-Kernel Methods**

▶ ~~§11.4 (+ §7.2) + §∞.2.2~~

~~Lecture 3~~ **The Principle of Sparsity (Recurring)**

▶ ~~§4, §8, §11.2, §∞.3~~

~~Lecture 4~~ **The Principle of Criticality**

▶ ~~§5, §9, §11.3, §∞.1, + §10.3.1~~

~~Lecture 5~~ **The End of Training, & More**

▶ ~~§∞.2.3 + §10.3.2 + §A.3 + §ε~~

The End of Training

$$\begin{aligned}
 & \mathbf{z}_{i;\beta}(t = \infty) \\
 = & \mathbf{z}_{i;\beta} - \sum_{j,k,\tilde{\alpha}_1,\tilde{\alpha}_2} \hat{H}_{ij;\beta\tilde{\alpha}_1} (\hat{H}^{-1})_{jk}^{\tilde{\alpha}_1\tilde{\alpha}_2} (\mathbf{z}_{k;\tilde{\alpha}_2} - y_{k;\tilde{\alpha}_2}) \\
 & + \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4} \left[\widehat{dH}_{j_1j_2;\tilde{\alpha}_1\beta\tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5,\tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{j_1j_2;\tilde{\alpha}_1\tilde{\alpha}_6\tilde{\alpha}_2} \right] Z_A^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4} (\mathbf{z}_{j_1;\tilde{\alpha}_3} - y_{j_1;\tilde{\alpha}_3}) (\mathbf{z}_{j_2;\tilde{\alpha}_4} - y_{j_2;\tilde{\alpha}_4}) \\
 & + \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4} \left[\widehat{dH}_{ij_1j_2;\beta\tilde{\alpha}_1\tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5,\tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{ij_1j_2;\tilde{\alpha}_6\tilde{\alpha}_1\tilde{\alpha}_2} \right] Z_B^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4} (\mathbf{z}_{j_1;\tilde{\alpha}_3} - y_{j_1;\tilde{\alpha}_3}) (\mathbf{z}_{j_2;\tilde{\alpha}_4} - y_{j_2;\tilde{\alpha}_4}) \\
 & + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[\widehat{dd}_I H_{j_1j_2j_3;\tilde{\alpha}_1\beta\tilde{\alpha}_2\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_I H_{j_1j_2j_3;\tilde{\alpha}_1\tilde{\alpha}_8\tilde{\alpha}_2\tilde{\alpha}_3} \right] Z_{IA}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (\mathbf{z}_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (\mathbf{z}_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (\mathbf{z}_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[\widehat{dd}_I H_{ij_1j_2j_3;\beta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_I H_{ij_1j_2j_3;\tilde{\alpha}_8\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} \right] Z_{IB}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (\mathbf{z}_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (\mathbf{z}_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (\mathbf{z}_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[\widehat{dd}_{II} H_{j_1j_2j_3;\tilde{\alpha}_1\tilde{\alpha}_2\beta\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{II} H_{j_1j_2j_3;\tilde{\alpha}_1\tilde{\alpha}_8\tilde{\alpha}_2\tilde{\alpha}_3} \right] Z_{IIA}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (\mathbf{z}_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (\mathbf{z}_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (\mathbf{z}_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[\widehat{dd}_{II} H_{ij_1j_2j_3;\beta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{II} H_{ij_1j_2j_3;\tilde{\alpha}_8\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} \right] Z_{IIB}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (\mathbf{z}_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (\mathbf{z}_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (\mathbf{z}_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & + \mathcal{O}\left(\frac{1}{n^2}\right)
 \end{aligned}$$

The End of Training

$$\begin{aligned}
 & \mathbf{z}_{i;\beta}(t = \infty) \\
 = & \mathbf{z}_{i;\beta} - \sum_{j,k,\tilde{\alpha}_1,\tilde{\alpha}_2} \hat{H}_{ij;\beta\tilde{\alpha}_1} (\hat{H}^{-1})_{jk}^{\tilde{\alpha}_1\tilde{\alpha}_2} (\mathbf{z}_{k;\tilde{\alpha}_2} - y_{k;\tilde{\alpha}_2}) \\
 & + \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4} \left[\widehat{dH}_{j_1j_2;\tilde{\alpha}_1\beta\tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5,\tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{j_1j_2;\tilde{\alpha}_1\tilde{\alpha}_6\tilde{\alpha}_2} \right] Z_A^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4} (\mathbf{z}_{j_1;\tilde{\alpha}_3} - y_{j_1;\tilde{\alpha}_3}) (\mathbf{z}_{j_2;\tilde{\alpha}_4} - y_{j_2;\tilde{\alpha}_4}) \\
 & + \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4} \left[\widehat{dH}_{ij_1j_2;\beta\tilde{\alpha}_1\tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5,\tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{ij_1j_2;\tilde{\alpha}_6\tilde{\alpha}_1\tilde{\alpha}_2} \right] Z_B^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4} (\mathbf{z}_{j_1;\tilde{\alpha}_3} - y_{j_1;\tilde{\alpha}_3}) (\mathbf{z}_{j_2;\tilde{\alpha}_4} - y_{j_2;\tilde{\alpha}_4}) \\
 & + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[\widehat{dd}_I H_{j_1j_2j_3;\tilde{\alpha}_1\beta\tilde{\alpha}_2\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_I H_{j_1j_2j_3;\tilde{\alpha}_1\tilde{\alpha}_8\tilde{\alpha}_2\tilde{\alpha}_3} \right] Z_{IA}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (\mathbf{z}_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (\mathbf{z}_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (\mathbf{z}_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[\widehat{dd}_I H_{ij_1j_2j_3;\beta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_I H_{ij_1j_2j_3;\tilde{\alpha}_8\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} \right] Z_{IB}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (\mathbf{z}_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (\mathbf{z}_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (\mathbf{z}_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[\widehat{dd}_{II} H_{j_1j_2j_3;\tilde{\alpha}_1\tilde{\alpha}_2\beta\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{II} H_{j_1j_2j_3;\tilde{\alpha}_1\tilde{\alpha}_8\tilde{\alpha}_2\tilde{\alpha}_3} \right] Z_{IIA}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (\mathbf{z}_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (\mathbf{z}_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (\mathbf{z}_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[\widehat{dd}_{II} H_{ij_1j_2j_3;\beta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{II} H_{ij_1j_2j_3;\tilde{\alpha}_8\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} \right] Z_{IIB}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (\mathbf{z}_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (\mathbf{z}_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (\mathbf{z}_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & + \mathcal{O}\left(\frac{1}{n^2}\right)
 \end{aligned}$$

Complete Inverse of the Stochastic NTK

Here, we needed the complete inverse of the stochastic NTK:

$$\sum_{j, \tilde{\alpha}_2} \left(\hat{H}^{-1} \right)_{ij}^{\tilde{\alpha}_1 \tilde{\alpha}_2} \hat{H}_{jk; \tilde{\alpha}_2 \tilde{\alpha}_3} = \delta_{ik} \delta_{\tilde{\alpha}_3}^{\tilde{\alpha}_1} .$$

Complete Inverse of the Stochastic NTK

Here, we needed the complete inverse of the stochastic NTK:

$$\sum_{j, \tilde{\alpha}_2} \left(\widehat{H}^{-1} \right)_{ij}^{\tilde{\alpha}_1 \tilde{\alpha}_2} \widehat{H}_{jk; \tilde{\alpha}_2 \tilde{\alpha}_3} = \delta_{ik} \delta_{\tilde{\alpha}_3}^{\tilde{\alpha}_1} .$$

Defining a mean and fluctuation as

$$\widehat{H}_{i_1 i_2; \alpha_1 \alpha_2} \equiv \delta_{i_1 i_2} H_{\alpha_1 \alpha_2} + \widehat{\Delta H}_{i_1 i_2; \alpha_1 \alpha_2} ,$$

we can then expand around the fluctuation to get:

$$\begin{aligned} \left(\widehat{H}^{-1} \right)_{ij}^{\tilde{\alpha}_1 \tilde{\alpha}_2} &= \delta_{ij} \widetilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} \widetilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_3} \widehat{\Delta H}_{ij; \tilde{\alpha}_3 \tilde{\alpha}_4} \widetilde{H}^{\tilde{\alpha}_4 \tilde{\alpha}_2} \\ &+ \sum_{k=1}^{n_L} \sum_{\tilde{\alpha}_3, \dots, \tilde{\alpha}_6 \in \mathcal{A}} \widetilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_3} \widehat{\Delta H}_{ik; \tilde{\alpha}_3 \tilde{\alpha}_4} \widetilde{H}^{\tilde{\alpha}_4 \tilde{\alpha}_5} \widehat{\Delta H}_{kj; \tilde{\alpha}_5 \tilde{\alpha}_6} \widetilde{H}^{\tilde{\alpha}_6 \tilde{\alpha}_2} \\ &+ O(\Delta^3) . \end{aligned}$$

The End of Training

$$\begin{aligned}
 & z_{i,\beta}(t = \infty) \\
 = & z_{i,\beta} - \sum_{j,k,\tilde{\alpha}_1,\tilde{\alpha}_2} \hat{H}_{ij;\beta\tilde{\alpha}_1} (\hat{H}^{-1})_{jk}^{\tilde{\alpha}_1\tilde{\alpha}_2} (z_{k;\tilde{\alpha}_2} - y_{k;\tilde{\alpha}_2}) \\
 & + \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4} \left[\widehat{dH}_{j_1j_2;\tilde{\alpha}_1\beta\tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5,\tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{j_1j_2;\tilde{\alpha}_1\tilde{\alpha}_6\tilde{\alpha}_2} \right] Z_A^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4} (z_{j_1;\tilde{\alpha}_3} - y_{j_1;\tilde{\alpha}_3}) (z_{j_2;\tilde{\alpha}_4} - y_{j_2;\tilde{\alpha}_4}) \\
 & + \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4} \left[\widehat{dH}_{j_1j_2;\beta\tilde{\alpha}_1\tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5,\tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{j_1j_2;\tilde{\alpha}_6\tilde{\alpha}_1\tilde{\alpha}_2} \right] Z_B^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4} (z_{j_1;\tilde{\alpha}_3} - y_{j_1;\tilde{\alpha}_3}) (z_{j_2;\tilde{\alpha}_4} - y_{j_2;\tilde{\alpha}_4}) \\
 & + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[\widehat{dd}_1 H_{j_1j_2j_3;\tilde{\alpha}_1\beta\tilde{\alpha}_2\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_1 H_{j_1j_2j_3;\tilde{\alpha}_1\tilde{\alpha}_8\tilde{\alpha}_2\tilde{\alpha}_3} \right] Z_{IA}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (z_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (z_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (z_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[\widehat{dd}_1 H_{j_1j_2j_3;\beta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_1 H_{j_1j_2j_3;\tilde{\alpha}_8\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} \right] Z_{IB}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (z_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (z_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (z_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[\widehat{dd}_1 H_{j_1j_2j_3;\tilde{\alpha}_1\tilde{\alpha}_2\beta\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_1 H_{j_1j_2j_3;\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_8\tilde{\alpha}_3} \right] Z_{IIA}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (z_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (z_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (z_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[\widehat{dd}_1 H_{j_1j_2j_3;\beta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7,\tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_1 H_{j_1j_2j_3;\tilde{\alpha}_8\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} \right] Z_{IIB}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (z_{j_1;\tilde{\alpha}_4} - y_{j_1;\tilde{\alpha}_4}) (z_{j_2;\tilde{\alpha}_5} - y_{j_2;\tilde{\alpha}_5}) (z_{j_3;\tilde{\alpha}_6} - y_{j_3;\tilde{\alpha}_6}) \\
 & + O\left(\frac{1}{n^2}\right)
 \end{aligned}$$

The End of Training

$$\begin{aligned}
 & \mathbf{z}_{i;\beta}(t = \infty) \\
 = & \mathbf{z}_{i;\beta} - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\beta\tilde{\alpha}_1} \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} (\mathbf{z}_{i;\tilde{\alpha}_2} - \mathbf{y}_{i;\tilde{\alpha}_2}) \\
 & + \sum_{j, \tilde{\alpha}_1, \tilde{\alpha}_2} \left[\widehat{\Delta H}_{ij;\beta\tilde{\alpha}_1} - \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} H_{\beta\tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4} \widehat{\Delta H}_{ij;\tilde{\alpha}_4\tilde{\alpha}_1} \right] \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} (\mathbf{z}_{j;\tilde{\alpha}_2} - \mathbf{y}_{j;\tilde{\alpha}_2}) \\
 & - \sum_{\substack{j,k \\ \tilde{\alpha}_1, \dots, \tilde{\alpha}_4}} \left[\widehat{\Delta H}_{ij;\beta\tilde{\alpha}_1} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{\Delta H}_{ij;\tilde{\alpha}_6\tilde{\alpha}_1} \right] \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} \widehat{\Delta H}_{jk;\tilde{\alpha}_2\tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4} (\mathbf{z}_{k;\tilde{\alpha}_4} - \mathbf{y}_{k;\tilde{\alpha}_4}) \\
 & + \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \left[\widehat{dH}_{j_1 j_2; \tilde{\alpha}_1 \beta \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{j_1 j_2; \tilde{\alpha}_1 \tilde{\alpha}_6 \tilde{\alpha}_2} \right] Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (\mathbf{z}_{j_1; \tilde{\alpha}_3} - \mathbf{y}_{j_1; \tilde{\alpha}_3}) (\mathbf{z}_{j_2; \tilde{\alpha}_4} - \mathbf{y}_{j_2; \tilde{\alpha}_4}) \\
 & + \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \left[\widehat{dH}_{ij_1 j_2; \beta \tilde{\alpha}_1 \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{ij_1 j_2; \tilde{\alpha}_6 \tilde{\alpha}_1 \tilde{\alpha}_2} \right] Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (\mathbf{z}_{j_1; \tilde{\alpha}_3} - \mathbf{y}_{j_1; \tilde{\alpha}_3}) (\mathbf{z}_{j_2; \tilde{\alpha}_4} - \mathbf{y}_{j_2; \tilde{\alpha}_4}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_1 H_{j_1 j_2 j_3; \tilde{\alpha}_1 \beta \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_1 H_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_8 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_1 H_{ij_1 j_2 j_3; \beta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_1 H_{ij_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{II} H_{j_1 j_2 ij_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \beta \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{II} H_{j_1 j_2 ij_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_8 \tilde{\alpha}_3} \right] Z_{IIA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{II} H_{ij_1 j_2 j_3; \beta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{II} H_{ij_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IIB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \mathcal{O}\left(\frac{1}{n^2}\right)
 \end{aligned}$$

The End of Training

$$\begin{aligned}
 & z_{i;\beta}(t = \infty) \\
 = & z_{i;\beta} - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\beta\tilde{\alpha}_1} \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2}(z_{i;\tilde{\alpha}_2} - y_{i;\tilde{\alpha}_2}) \\
 & + \sum_{j, \tilde{\alpha}_1, \tilde{\alpha}_2} \left[\widehat{\Delta H}_{ij;\beta\tilde{\alpha}_1} - \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} H_{\beta\tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4} \widehat{\Delta H}_{ij;\tilde{\alpha}_4\tilde{\alpha}_1} \right] \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2}(z_{j;\tilde{\alpha}_2} - y_{j;\tilde{\alpha}_2}) \\
 & - \sum_{\substack{j,k \\ \tilde{\alpha}_1, \dots, \tilde{\alpha}_4}} \left[\widehat{\Delta H}_{ij;\beta\tilde{\alpha}_1} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{\Delta H}_{ij;\tilde{\alpha}_6\tilde{\alpha}_1} \right] \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} \widehat{\Delta H}_{jk;\tilde{\alpha}_2\tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4} (z_{k;\tilde{\alpha}_4} - y_{k;\tilde{\alpha}_4}) \\
 & + \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \left[\widehat{dH}_{j_1 j_2; \tilde{\alpha}_1 \beta \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{j_1 j_2; \tilde{\alpha}_1 \tilde{\alpha}_6 \tilde{\alpha}_2} \right] Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (z_{j_1; \tilde{\alpha}_3} - y_{j_1; \tilde{\alpha}_3}) (z_{j_2; \tilde{\alpha}_4} - y_{j_2; \tilde{\alpha}_4}) \\
 & + \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \left[\widehat{dH}_{ij_1 j_2; \beta \tilde{\alpha}_1 \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{ij_1 j_2; \tilde{\alpha}_6 \tilde{\alpha}_1 \tilde{\alpha}_2} \right] Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (z_{j_1; \tilde{\alpha}_3} - y_{j_1; \tilde{\alpha}_3}) (z_{j_2; \tilde{\alpha}_4} - y_{j_2; \tilde{\alpha}_4}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{j_1 j_2 j_3; \tilde{\alpha}_1 \beta \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_8 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4}) (z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5}) (z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{ij_1 j_2 j_3; \beta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{ij_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4}) (z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5}) (z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \beta \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_8 \tilde{\alpha}_3} \right] Z_{IIA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4}) (z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5}) (z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{ij_1 j_2 j_3; \beta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{ij_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IIB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4}) (z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5}) (z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}) \\
 & + \mathcal{O}\left(\frac{1}{n^2}\right)
 \end{aligned}$$

Two Extra Differentials??

$$\widehat{dd_I H}_{i_0 i_1 i_2 i_3; \delta_0 \delta_1 \delta_2 \delta_3}^{(\ell)}$$

$$\equiv \sum_{\ell_1, \ell_2, \ell_3=1}^{\ell} \sum_{\substack{\mu_1, \nu_1, \\ \mu_2, \nu_2, \\ \mu_3, \nu_3}} \lambda_{\mu_1 \nu_1}^{(\ell_1)} \lambda_{\mu_2 \nu_2}^{(\ell_2)} \lambda_{\mu_3 \nu_3}^{(\ell_3)} \frac{d^3 z_{i_0; \delta_0}^{(\ell)}}{d\theta_{\mu_1}^{(\ell_1)} d\theta_{\mu_2}^{(\ell_2)} d\theta_{\mu_3}^{(\ell_3)}} \frac{dz_{i_1; \delta_1}^{(\ell)}}{d\theta_{\nu_1}^{(\ell_1)}} \frac{dz_{i_2; \delta_2}^{(\ell)}}{d\theta_{\nu_2}^{(\ell_2)}} \frac{dz_{i_3; \delta_3}^{(\ell)}}{d\theta_{\nu_3}^{(\ell_3)}}$$

$$\widehat{dd_{II} H}_{i_1 i_2 i_3 i_4; \delta_1 \delta_2 \delta_3 \delta_4}^{(\ell)}$$

$$\equiv \sum_{\ell_1, \ell_2, \ell_3=1}^{\ell} \sum_{\substack{\mu_1, \nu_1, \\ \mu_2, \nu_2, \\ \mu_3, \nu_3}} \lambda_{\mu_1 \nu_1}^{(\ell_1)} \lambda_{\mu_2 \nu_2}^{(\ell_2)} \lambda_{\mu_3 \nu_3}^{(\ell_3)} \frac{d^2 z_{i_1; \delta_1}^{(\ell)}}{d\theta_{\mu_1}^{(\ell_1)} d\theta_{\mu_3}^{(\ell_3)}} \frac{d^2 z_{i_2; \delta_2}^{(\ell)}}{d\theta_{\mu_2}^{(\ell_2)} d\theta_{\nu_3}^{(\ell_3)}} \frac{dz_{i_3; \delta_3}^{(\ell)}}{d\theta_{\nu_1}^{(\ell_1)}} \frac{dz_{i_4; \delta_4}^{(\ell)}}{d\theta_{\nu_2}^{(\ell_2)}}$$

The End of Training

$$\begin{aligned}
 & z_{i;\beta}(t = \infty) \\
 = & z_{i;\beta} - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\beta\tilde{\alpha}_1} \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2}(z_{i;\tilde{\alpha}_2} - y_{i;\tilde{\alpha}_2}) \\
 & + \sum_{j, \tilde{\alpha}_1, \tilde{\alpha}_2} \left[\widehat{\Delta H}_{ij;\beta\tilde{\alpha}_1} - \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} H_{\beta\tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4} \widehat{\Delta H}_{ij;\tilde{\alpha}_4\tilde{\alpha}_1} \right] \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2}(z_{j;\tilde{\alpha}_2} - y_{j;\tilde{\alpha}_2}) \\
 & - \sum_{\substack{j,k \\ \tilde{\alpha}_1, \dots, \tilde{\alpha}_4}} \left[\widehat{\Delta H}_{ij;\beta\tilde{\alpha}_1} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{\Delta H}_{ij;\tilde{\alpha}_6\tilde{\alpha}_1} \right] \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} \widehat{\Delta H}_{jk;\tilde{\alpha}_2\tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4}(z_{k;\tilde{\alpha}_4} - y_{k;\tilde{\alpha}_4}) \\
 & + \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \left[\widehat{dH}_{j_1 j_2; \tilde{\alpha}_1 \beta \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{j_1 j_2; \tilde{\alpha}_1 \tilde{\alpha}_6 \tilde{\alpha}_2} \right] Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4}(z_{j_1; \tilde{\alpha}_3} - y_{j_1; \tilde{\alpha}_3})(z_{j_2; \tilde{\alpha}_4} - y_{j_2; \tilde{\alpha}_4}) \\
 & + \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \left[\widehat{dH}_{ij_1 j_2; \beta \tilde{\alpha}_1 \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{ij_1 j_2; \tilde{\alpha}_6 \tilde{\alpha}_1 \tilde{\alpha}_2} \right] Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4}(z_{j_1; \tilde{\alpha}_3} - y_{j_1; \tilde{\alpha}_3})(z_{j_2; \tilde{\alpha}_4} - y_{j_2; \tilde{\alpha}_4}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_1 H_{j_1 j_2 j_3; \tilde{\alpha}_1 \beta \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_1 H_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_8 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}(z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4})(z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5})(z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_1 H_{ij_1 j_2 j_3; \beta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_1 H_{ij_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}(z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4})(z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5})(z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{II} H_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \beta \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{II} H_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_8 \tilde{\alpha}_3} \right] Z_{IIA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}(z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4})(z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5})(z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{II} H_{ij_1 j_2 j_3; \beta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{II} H_{ij_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IIB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}(z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4})(z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5})(z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}) \\
 & + \mathcal{O}\left(\frac{1}{n^2}\right)
 \end{aligned}$$

The End of Training

$$\begin{aligned}
 & \mathbf{z}_{i;\beta}(t = \infty) \\
 = & \mathbf{z}_{i;\beta} - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\beta\tilde{\alpha}_1} \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2}(\mathbf{z}_{i;\tilde{\alpha}_2} - \mathbf{y}_{i;\tilde{\alpha}_2}) \\
 & + \sum_{j, \tilde{\alpha}_1, \tilde{\alpha}_2} \left[\widehat{\Delta H}_{ij;\beta\tilde{\alpha}_1} - \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} H_{\beta\tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4} \widehat{\Delta H}_{ij;\tilde{\alpha}_4\tilde{\alpha}_1} \right] \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2}(\mathbf{z}_{j;\tilde{\alpha}_2} - \mathbf{y}_{j;\tilde{\alpha}_2}) \\
 & - \sum_{\substack{j,k \\ \tilde{\alpha}_1, \dots, \tilde{\alpha}_4}} \left[\widehat{\Delta H}_{ij;\beta\tilde{\alpha}_1} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{\Delta H}_{ij;\tilde{\alpha}_6\tilde{\alpha}_1} \right] \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} \widehat{\Delta H}_{jk;\tilde{\alpha}_2\tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4}(\mathbf{z}_{k;\tilde{\alpha}_4} - \mathbf{y}_{k;\tilde{\alpha}_4}) \\
 & + \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \left[\widehat{dH}_{j_1 j_2; \tilde{\alpha}_1 \beta \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{j_1 j_2; \tilde{\alpha}_1 \tilde{\alpha}_6 \tilde{\alpha}_2} \right] Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4}(\mathbf{z}_{j_1; \tilde{\alpha}_3} - \mathbf{y}_{j_1; \tilde{\alpha}_3})(\mathbf{z}_{j_2; \tilde{\alpha}_4} - \mathbf{y}_{j_2; \tilde{\alpha}_4}) \\
 & + \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \left[\widehat{dH}_{ij_1 j_2; \beta \tilde{\alpha}_1 \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{ij_1 j_2; \tilde{\alpha}_6 \tilde{\alpha}_1 \tilde{\alpha}_2} \right] Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4}(\mathbf{z}_{j_1; \tilde{\alpha}_3} - \mathbf{y}_{j_1; \tilde{\alpha}_3})(\mathbf{z}_{j_2; \tilde{\alpha}_4} - \mathbf{y}_{j_2; \tilde{\alpha}_4}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_1 H_{j_1 j_2 j_3; \tilde{\alpha}_1 \beta \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_1 H_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_8 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}(\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4})(\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5})(\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_1 H_{ij_1 j_2 j_3; \beta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_1 H_{ij_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}(\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4})(\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5})(\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{II} H_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \beta \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{II} H_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_8 \tilde{\alpha}_3} \right] Z_{IIA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}(\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4})(\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5})(\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{II} H_{ij_1 j_2 j_3; \beta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{II} H_{ij_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IIB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}(\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4})(\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5})(\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \mathcal{O}\left(\frac{1}{n^2}\right)
 \end{aligned}$$

The End of Training

$$\begin{aligned}
 & \mathbf{z}_{i;\beta}(t = \infty) \\
 = & \mathbf{z}_{i;\beta} - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\beta\tilde{\alpha}_1} \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} (\mathbf{z}_{i;\tilde{\alpha}_2} - \mathbf{y}_{i;\tilde{\alpha}_2}) \\
 & + \sum_{j, \tilde{\alpha}_1, \tilde{\alpha}_2} \left[\widehat{\Delta H}_{ij;\beta\tilde{\alpha}_1} - \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} H_{\beta\tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4} \widehat{\Delta H}_{ij;\tilde{\alpha}_4\tilde{\alpha}_1} \right] \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} (\mathbf{z}_{j;\tilde{\alpha}_2} - \mathbf{y}_{j;\tilde{\alpha}_2}) \\
 & - \sum_{\substack{j,k \\ \tilde{\alpha}_1, \dots, \tilde{\alpha}_4}} \left[\widehat{\Delta H}_{ij;\beta\tilde{\alpha}_1} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{\Delta H}_{ij;\tilde{\alpha}_6\tilde{\alpha}_1} \right] \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} \widehat{\Delta H}_{jk;\tilde{\alpha}_2\tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4} (\mathbf{z}_{k;\tilde{\alpha}_4} - \mathbf{y}_{k;\tilde{\alpha}_4}) \\
 & + \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \left[\widehat{dH}_{j_1 j_2; \tilde{\alpha}_1 \beta \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{j_1 j_2; \tilde{\alpha}_1 \tilde{\alpha}_6 \tilde{\alpha}_2} \right] Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (\mathbf{z}_{j_1; \tilde{\alpha}_3} - \mathbf{y}_{j_1; \tilde{\alpha}_3}) (\mathbf{z}_{j_2; \tilde{\alpha}_4} - \mathbf{y}_{j_2; \tilde{\alpha}_4}) \\
 & + \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \left[\widehat{dH}_{ij_1 j_2; \beta \tilde{\alpha}_1 \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{ij_1 j_2; \tilde{\alpha}_6 \tilde{\alpha}_1 \tilde{\alpha}_2} \right] Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (\mathbf{z}_{j_1; \tilde{\alpha}_3} - \mathbf{y}_{j_1; \tilde{\alpha}_3}) (\mathbf{z}_{j_2; \tilde{\alpha}_4} - \mathbf{y}_{j_2; \tilde{\alpha}_4}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{j_1 j_2 j_3; \tilde{\alpha}_1 \beta \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_8 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{ij_1 j_2 j_3; \beta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{ij_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \beta \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_8 \tilde{\alpha}_3} \right] Z_{IIA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{ij_1 j_2 j_3; \beta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{ij_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IIB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \mathcal{O}\left(\frac{1}{n^2}\right)
 \end{aligned}$$

The End of Training

$$\begin{aligned}
 & \mathbf{z}_{i;\beta}(t = \infty) \\
 = & \mathbf{z}_{i;\beta} - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\beta\tilde{\alpha}_1} \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} (\mathbf{z}_{i;\tilde{\alpha}_2} - \mathbf{y}_{i;\tilde{\alpha}_2}) \\
 & + \sum_{j, \tilde{\alpha}_1, \tilde{\alpha}_2} \left[\widehat{\Delta H}_{ij;\beta\tilde{\alpha}_1} - \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} H_{\beta\tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4} \widehat{\Delta H}_{ij;\tilde{\alpha}_4\tilde{\alpha}_1} \right] \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} (\mathbf{z}_{j;\tilde{\alpha}_2} - \mathbf{y}_{j;\tilde{\alpha}_2}) \\
 & - \sum_{\substack{j,k \\ \tilde{\alpha}_1, \dots, \tilde{\alpha}_4}} \left[\widehat{\Delta H}_{ij;\beta\tilde{\alpha}_1} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{\Delta H}_{ij;\tilde{\alpha}_6\tilde{\alpha}_1} \right] \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} \widehat{\Delta H}_{jk;\tilde{\alpha}_2\tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4} (\mathbf{z}_{k;\tilde{\alpha}_4} - \mathbf{y}_{k;\tilde{\alpha}_4}) \\
 & + \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \left[\widehat{dH}_{j_1 j_2; \tilde{\alpha}_1 \beta \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{j_1 j_2; \tilde{\alpha}_1 \tilde{\alpha}_6 \tilde{\alpha}_2} \right] Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (\mathbf{z}_{j_1; \tilde{\alpha}_3} - \mathbf{y}_{j_1; \tilde{\alpha}_3}) (\mathbf{z}_{j_2; \tilde{\alpha}_4} - \mathbf{y}_{j_2; \tilde{\alpha}_4}) \\
 & + \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \left[\widehat{dH}_{ij_1 j_2; \beta \tilde{\alpha}_1 \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{ij_1 j_2; \tilde{\alpha}_6 \tilde{\alpha}_1 \tilde{\alpha}_2} \right] Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (\mathbf{z}_{j_1; \tilde{\alpha}_3} - \mathbf{y}_{j_1; \tilde{\alpha}_3}) (\mathbf{z}_{j_2; \tilde{\alpha}_4} - \mathbf{y}_{j_2; \tilde{\alpha}_4}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{ddH}_{j_1 j_2 j_3; \tilde{\alpha}_1 \beta \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{ddH}_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_8 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{ddH}_{ij_1 j_2 j_3; \beta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{ddH}_{ij_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{ddH}_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \beta \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{ddH}_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_8 \tilde{\alpha}_3} \right] Z_{IIA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{ddH}_{ij_1 j_2 j_3; \beta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{ddH}_{ij_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IIB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \mathcal{O}\left(\frac{1}{n^2}\right)
 \end{aligned}$$

The End of Training

$$\begin{aligned}
 & \mathbf{z}_{i;\beta}(t = \infty) \\
 = & \mathbf{z}_{i;\beta} - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\beta\tilde{\alpha}_1} \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} (\mathbf{z}_{i;\tilde{\alpha}_2} - \mathbf{y}_{i;\tilde{\alpha}_2}) \\
 & + \sum_{j, \tilde{\alpha}_1, \tilde{\alpha}_2} \left[\widehat{\Delta H}_{ij;\beta\tilde{\alpha}_1} - \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} H_{\beta\tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4} \widehat{\Delta H}_{ij;\tilde{\alpha}_4\tilde{\alpha}_1} \right] \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} (\mathbf{z}_{j;\tilde{\alpha}_2} - \mathbf{y}_{j;\tilde{\alpha}_2}) \\
 & - \sum_{\substack{j,k \\ \tilde{\alpha}_1, \dots, \tilde{\alpha}_4}} \left[\widehat{\Delta H}_{ij;\beta\tilde{\alpha}_1} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{\Delta H}_{ij;\tilde{\alpha}_6\tilde{\alpha}_1} \right] \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} \widehat{\Delta H}_{jk;\tilde{\alpha}_2\tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4} (\mathbf{z}_{k;\tilde{\alpha}_4} - \mathbf{y}_{k;\tilde{\alpha}_4}) \\
 & + \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \left[\widehat{dH}_{j_1 j_2; \tilde{\alpha}_1 \beta \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{j_1 j_2; \tilde{\alpha}_1 \tilde{\alpha}_6 \tilde{\alpha}_2} \right] Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (\mathbf{z}_{j_1; \tilde{\alpha}_3} - \mathbf{y}_{j_1; \tilde{\alpha}_3}) (\mathbf{z}_{j_2; \tilde{\alpha}_4} - \mathbf{y}_{j_2; \tilde{\alpha}_4}) \\
 & + \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \left[\widehat{dH}_{ij_1 j_2; \beta \tilde{\alpha}_1 \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{ij_1 j_2; \tilde{\alpha}_6 \tilde{\alpha}_1 \tilde{\alpha}_2} \right] Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (\mathbf{z}_{j_1; \tilde{\alpha}_3} - \mathbf{y}_{j_1; \tilde{\alpha}_3}) (\mathbf{z}_{j_2; \tilde{\alpha}_4} - \mathbf{y}_{j_2; \tilde{\alpha}_4}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_I H_{j_1 j_2 j_3; \tilde{\alpha}_1 \beta \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_I H_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_8 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_I H_{ij_1 j_2 j_3; \beta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_I H_{ij_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{II} H_{j_1 j_2 ij_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \beta \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{II} H_{j_1 j_2 ij_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_8 \tilde{\alpha}_3} \right] Z_{IIA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{II} H_{ij_1 j_2 j_3; \beta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{II} H_{ij_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IIB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \mathcal{O}\left(\frac{1}{n^2}\right)
 \end{aligned}$$

The End of Training

$$\begin{aligned}
 & \mathbf{z}_{i;\beta}(t = \infty) \\
 = & \mathbf{z}_{i;\beta} - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\beta\tilde{\alpha}_1} \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} (\mathbf{z}_{i;\tilde{\alpha}_2} - \mathbf{y}_{i;\tilde{\alpha}_2}) \\
 & + \sum_{j, \tilde{\alpha}_1, \tilde{\alpha}_2} \left[\widehat{\Delta H}_{ij;\beta\tilde{\alpha}_1} - \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} H_{\beta\tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4} \widehat{\Delta H}_{ij;\tilde{\alpha}_4\tilde{\alpha}_1} \right] \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} (\mathbf{z}_{j;\tilde{\alpha}_2} - \mathbf{y}_{j;\tilde{\alpha}_2}) \\
 & - \sum_{\substack{j,k \\ \tilde{\alpha}_1, \dots, \tilde{\alpha}_4}} \left[\widehat{\Delta H}_{ij;\beta\tilde{\alpha}_1} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{\Delta H}_{ij;\tilde{\alpha}_6\tilde{\alpha}_1} \right] \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} \widehat{\Delta H}_{jk;\tilde{\alpha}_2\tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4} (\mathbf{z}_{k;\tilde{\alpha}_4} - \mathbf{y}_{k;\tilde{\alpha}_4}) \\
 & + \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \left[\widehat{dH}_{j_1 j_2; \tilde{\alpha}_1 \beta \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{j_1 j_2; \tilde{\alpha}_1 \tilde{\alpha}_6 \tilde{\alpha}_2} \right] Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (\mathbf{z}_{j_1; \tilde{\alpha}_3} - \mathbf{y}_{j_1; \tilde{\alpha}_3}) (\mathbf{z}_{j_2; \tilde{\alpha}_4} - \mathbf{y}_{j_2; \tilde{\alpha}_4}) \\
 & + \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \left[\widehat{dH}_{ij_1 j_2; \beta \tilde{\alpha}_1 \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{ij_1 j_2; \tilde{\alpha}_6 \tilde{\alpha}_1 \tilde{\alpha}_2} \right] Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (\mathbf{z}_{j_1; \tilde{\alpha}_3} - \mathbf{y}_{j_1; \tilde{\alpha}_3}) (\mathbf{z}_{j_2; \tilde{\alpha}_4} - \mathbf{y}_{j_2; \tilde{\alpha}_4}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{j_1 j_2 j_3; \tilde{\alpha}_1 \beta \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_8 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{j_1 j_2 j_3; \beta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{j_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \beta \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_8 \tilde{\alpha}_3} \right] Z_{IIA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{j_1 j_2 j_3; \beta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{j_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IIB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \mathcal{O}\left(\frac{1}{n^2}\right)
 \end{aligned}$$

The End of Training

$$\begin{aligned}
 & \mathbf{z}_{i;\beta}(t = \infty) \\
 = & \mathbf{z}_{i;\beta} - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\beta\tilde{\alpha}_1} \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} (\mathbf{z}_{i;\tilde{\alpha}_2} - \mathbf{y}_{i;\tilde{\alpha}_2}) \\
 & + \sum_{j, \tilde{\alpha}_1, \tilde{\alpha}_2} \left[\widehat{\Delta H}_{ij;\beta\tilde{\alpha}_1} - \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} H_{\beta\tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4} \widehat{\Delta H}_{ij;\tilde{\alpha}_4\tilde{\alpha}_1} \right] \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} (\mathbf{z}_{j;\tilde{\alpha}_2} - \mathbf{y}_{j;\tilde{\alpha}_2}) \\
 & - \sum_{\substack{j,k \\ \tilde{\alpha}_1, \dots, \tilde{\alpha}_4}} \left[\widehat{\Delta H}_{ij;\beta\tilde{\alpha}_1} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{\Delta H}_{ij;\tilde{\alpha}_6\tilde{\alpha}_1} \right] \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} \widehat{\Delta H}_{jk;\tilde{\alpha}_2\tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4} (\mathbf{z}_{k;\tilde{\alpha}_4} - \mathbf{y}_{k;\tilde{\alpha}_4}) \\
 & + \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \left[\widehat{dH}_{j_1 j_2; \tilde{\alpha}_1 \beta \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{j_1 j_2; \tilde{\alpha}_1 \tilde{\alpha}_6 \tilde{\alpha}_2} \right] Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (\mathbf{z}_{j_1; \tilde{\alpha}_3} - \mathbf{y}_{j_1; \tilde{\alpha}_3}) (\mathbf{z}_{j_2; \tilde{\alpha}_4} - \mathbf{y}_{j_2; \tilde{\alpha}_4}) \\
 & + \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \left[\widehat{dH}_{ij_1 j_2; \beta \tilde{\alpha}_1 \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\beta\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{dH}_{ij_1 j_2; \tilde{\alpha}_6 \tilde{\alpha}_1 \tilde{\alpha}_2} \right] Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (\mathbf{z}_{j_1; \tilde{\alpha}_3} - \mathbf{y}_{j_1; \tilde{\alpha}_3}) (\mathbf{z}_{j_2; \tilde{\alpha}_4} - \mathbf{y}_{j_2; \tilde{\alpha}_4}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_1 H_{j_1 j_2 j_3; \tilde{\alpha}_1 \beta \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_1 H_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_8 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_1 H_{ij_1 j_2 j_3; \beta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_1 H_{ij_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{II} H_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \beta \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{II} H_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_8 \tilde{\alpha}_3} \right] Z_{IIA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{dd}_{II} H_{ij_1 j_2 j_3; \beta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\beta\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{dd}_{II} H_{ij_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] Z_{IIB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (\mathbf{z}_{j_1; \tilde{\alpha}_4} - \mathbf{y}_{j_1; \tilde{\alpha}_4}) (\mathbf{z}_{j_2; \tilde{\alpha}_5} - \mathbf{y}_{j_2; \tilde{\alpha}_5}) (\mathbf{z}_{j_3; \tilde{\alpha}_6} - \mathbf{y}_{j_3; \tilde{\alpha}_6}) \\
 & + \mathcal{O}\left(\frac{1}{n^2}\right)
 \end{aligned}$$

The Algorithm Projectors

Gradient Flow (ODE limit):

$$Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \equiv Y_2^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4},$$

$$Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \equiv Y_2^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4},$$

$$Z_{IA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \equiv -Y_3^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6},$$

$$Z_{IB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \equiv -Y_3^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6},$$

$$Z_{IIA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \equiv -Y_3^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6},$$

$$Z_{IIB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \equiv -Y_3^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} - Y_3^{\tilde{\alpha}_2 \tilde{\alpha}_1 \tilde{\alpha}_3 \tilde{\alpha}_5 \tilde{\alpha}_4 \tilde{\alpha}_6} - Y_3^{\tilde{\alpha}_1 \tilde{\alpha}_3 \tilde{\alpha}_2 \tilde{\alpha}_4 \tilde{\alpha}_6 \tilde{\alpha}_5}.$$

Various Shorthands:

$$Y_2^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \equiv \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_2 \tilde{\alpha}_4} - \sum_{\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_2 \tilde{\alpha}_5} \chi_{II}^{\tilde{\alpha}_1 \tilde{\alpha}_5 \tilde{\alpha}_3 \tilde{\alpha}_4},$$

$$Y_3^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \equiv \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_4} \tilde{H}^{\tilde{\alpha}_2 \tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_3 \tilde{\alpha}_6} - \sum_{\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_2 \tilde{\alpha}_7} \chi_{II}^{\tilde{\alpha}_1 \tilde{\alpha}_7 \tilde{\alpha}_4 \tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_3 \tilde{\alpha}_6} - \sum_{\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_4} \tilde{H}^{\tilde{\alpha}_3 \tilde{\alpha}_7} \chi_{II}^{\tilde{\alpha}_2 \tilde{\alpha}_7 \tilde{\alpha}_5 \tilde{\alpha}_6} + \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8, \tilde{\alpha}_9} \tilde{H}^{\tilde{\alpha}_3 \tilde{\alpha}_9} \chi_{II}^{\tilde{\alpha}_2 \tilde{\alpha}_9 \tilde{\alpha}_7 \tilde{\alpha}_8} \chi_{III}^{\tilde{\alpha}_1 \tilde{\alpha}_7 \tilde{\alpha}_8 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6},$$

Inverting Tensors:

$$\delta_{\tilde{\alpha}_5}^{\tilde{\alpha}_1} \delta_{\tilde{\alpha}_6}^{\tilde{\alpha}_2} = \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} \chi_{II}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \left(\tilde{H}_{\tilde{\alpha}_3 \tilde{\alpha}_5} \delta_{\tilde{\alpha}_4 \tilde{\alpha}_6} + \delta_{\tilde{\alpha}_3 \tilde{\alpha}_5} \tilde{H}_{\tilde{\alpha}_4 \tilde{\alpha}_6} \right),$$

$$\delta_{\tilde{\alpha}_7}^{\tilde{\alpha}_1} \delta_{\tilde{\alpha}_8}^{\tilde{\alpha}_2} \delta_{\tilde{\alpha}_9}^{\tilde{\alpha}_3} = \sum_{\tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6} \chi_{III}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \left(\tilde{H}_{\tilde{\alpha}_4 \tilde{\alpha}_7} \delta_{\tilde{\alpha}_5 \tilde{\alpha}_8} \delta_{\tilde{\alpha}_6 \tilde{\alpha}_9} + \delta_{\tilde{\alpha}_4 \tilde{\alpha}_7} \tilde{H}_{\tilde{\alpha}_5 \tilde{\alpha}_8} \delta_{\tilde{\alpha}_6 \tilde{\alpha}_9} + \delta_{\tilde{\alpha}_4 \tilde{\alpha}_7} \delta_{\tilde{\alpha}_5 \tilde{\alpha}_8} \tilde{H}_{\tilde{\alpha}_6 \tilde{\alpha}_9} \right)$$

The Algorithm Projectors

Direct Optimization:

$$Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \equiv 0,$$

$$Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \equiv \frac{1}{2} \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_2 \tilde{\alpha}_4},$$

$$Z_{IA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \equiv 0,$$

$$Z_{IB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \equiv -\frac{1}{6} \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_4} \tilde{H}^{\tilde{\alpha}_2 \tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_3 \tilde{\alpha}_6}$$

$$Z_{IIA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \equiv 0,$$

$$Z_{IIB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \equiv 0.$$

Various Shorthands:

Inverting Tensors:

Mean Prediction

The *mean* prediction of a deep MLP is given by

$$\begin{aligned} \mathbb{E} \left[z_{i;\hat{\beta}}^{(L)}(\infty) \right] &\equiv m_{i;\hat{\beta}} \\ &= m_{i;\hat{\beta}}^{\text{NTK}} + \frac{1}{n_{L-1}} \left(m_{i;\hat{\beta}}^{\Delta\text{NTK}} + m_{i;\hat{\beta}}^{\text{dNTK}} + m_{i;\hat{\beta}}^{\text{ddNTK-I}} + m_{i;\hat{\beta}}^{\text{ddNTK-II}} \right) \\ &\quad - \frac{1}{n_{L-1}} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\hat{\beta}\tilde{\alpha}_1}^{(L)} \tilde{H}_{\tilde{\alpha}_1\tilde{\alpha}_2}^{(L)} \left(m_{i;\tilde{\alpha}_2}^{\Delta\text{NTK}} + m_{i;\tilde{\alpha}_2}^{\text{dNTK}} + m_{i;\tilde{\alpha}_2}^{\text{ddNTK-I}} + m_{i;\tilde{\alpha}_2}^{\text{ddNTK-II}} \right) \end{aligned}$$

Mean Prediction

The *mean* prediction of a deep MLP is given by

$$\begin{aligned} \mathbb{E} \left[z_{i;\dot{\beta}}^{(L)}(\infty) \right] &\equiv m_{i;\dot{\beta}} \\ &= m_{i;\dot{\beta}}^{\text{NTK}} + \frac{1}{n_{L-1}} \left(m_{i;\dot{\beta}}^{\Delta\text{NTK}} + m_{i;\dot{\beta}}^{\text{dNTK}} + m_{i;\dot{\beta}}^{\text{ddNTK-I}} + m_{i;\dot{\beta}}^{\text{ddNTK-II}} \right) \\ &\quad - \frac{1}{n_{L-1}} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\dot{\beta}\tilde{\alpha}_1}^{(L)} \tilde{H}_{\tilde{\alpha}_1\tilde{\alpha}_2}^{(L)} \left(m_{i;\tilde{\alpha}_2}^{\Delta\text{NTK}} + m_{i;\tilde{\alpha}_2}^{\text{dNTK}} + m_{i;\tilde{\alpha}_2}^{\text{ddNTK-I}} + m_{i;\tilde{\alpha}_2}^{\text{ddNTK-II}} \right) \end{aligned}$$

► The first term is the (neural tangent) *kernel prediction*:

$$m_{i;\dot{\beta}}^{\text{NTK}} \equiv \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\dot{\beta}\tilde{\alpha}_1}^{(L)} \tilde{H}_{\tilde{\alpha}_1\tilde{\alpha}_2}^{(L)} y_{i;\tilde{\alpha}_2}.$$

Mean Prediction

The *mean* prediction of a deep MLP is given by

$$\begin{aligned}\mathbb{E} \left[z_{i;\hat{\beta}}^{(L)}(\infty) \right] &\equiv m_{i;\hat{\beta}} \\ &= m_{i;\hat{\beta}}^{\text{NTK}} + \frac{1}{n_{L-1}} \left(m_{i;\hat{\beta}}^{\Delta\text{NTK}} + m_{i;\hat{\beta}}^{\text{dNTK}} + m_{i;\hat{\beta}}^{\text{ddNTK-I}} + m_{i;\hat{\beta}}^{\text{ddNTK-II}} \right) \\ &\quad - \frac{1}{n_{L-1}} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\hat{\beta}\tilde{\alpha}_1}^{(L)} \tilde{H}_{(\tilde{\alpha}_1\tilde{\alpha}_2)}^{(L)} \left(m_{i;\tilde{\alpha}_2}^{\Delta\text{NTK}} + m_{i;\tilde{\alpha}_2}^{\text{dNTK}} + m_{i;\tilde{\alpha}_2}^{\text{ddNTK-I}} + m_{i;\tilde{\alpha}_2}^{\text{ddNTK-II}} \right)\end{aligned}$$

- ▶ The first term is the (neural tangent) *kernel prediction*:

$$m_{i;\hat{\beta}}^{\text{NTK}} \equiv \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\hat{\beta}\tilde{\alpha}_1}^{(L)} \tilde{H}_{(\tilde{\alpha}_1\tilde{\alpha}_2)}^{(L)} y_{i;\tilde{\alpha}_2}.$$

- ▶ The **four other kinds of terms** give the $O(1/n)$ **corrections**.

Mean Prediction

The *mean* prediction of a deep MLP is given by

$$\begin{aligned} \mathbb{E} \left[z_{i;\dot{\beta}}^{(L)}(\infty) \right] &\equiv m_{i;\dot{\beta}} \\ &= m_{i;\dot{\beta}}^{\text{NTK}} + \frac{1}{n_{L-1}} \left(m_{i;\dot{\beta}}^{\Delta\text{NTK}} + m_{i;\dot{\beta}}^{\text{dNTK}} + m_{i;\dot{\beta}}^{\text{ddNTK-I}} + m_{i;\dot{\beta}}^{\text{ddNTK-II}} \right) \\ &\quad - \frac{1}{n_{L-1}} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\dot{\beta}\tilde{\alpha}_1}^{(L)} \tilde{H}_{\tilde{\alpha}_1\tilde{\alpha}_2}^{(L)} \left(m_{i;\tilde{\alpha}_2}^{\Delta\text{NTK}} + m_{i;\tilde{\alpha}_2}^{\text{dNTK}} + m_{i;\tilde{\alpha}_2}^{\text{ddNTK-I}} + m_{i;\tilde{\alpha}_2}^{\text{ddNTK-II}} \right) \end{aligned}$$

- ▶ Each network fits the training data with its own *particular* NTK, and so the resulting fully-trained particular output depends on the **NTK fluctuation**.

Mean Prediction

The *mean* prediction of a deep MLP is given by

$$\begin{aligned} \mathbb{E} \left[z_{i;\hat{\beta}}^{(L)}(\infty) \right] &\equiv m_{i;\hat{\beta}} \\ &= m_{i;\hat{\beta}}^{\text{NTK}} + \frac{1}{n_{L-1}} \left(m_{i;\hat{\beta}}^{\Delta\text{NTK}} + m_{i;\hat{\beta}}^{\text{dNTK}} + m_{i;\hat{\beta}}^{\text{ddNTK-I}} + m_{i;\hat{\beta}}^{\text{ddNTK-II}} \right) \\ &\quad - \frac{1}{n_{L-1}} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\tilde{\beta}\tilde{\alpha}_1}^{(L)} \tilde{H}_{\tilde{\alpha}_1\tilde{\alpha}_2}^{(L)} \left(m_{i;\tilde{\alpha}_2}^{\Delta\text{NTK}} + m_{i;\tilde{\alpha}_2}^{\text{dNTK}} + m_{i;\tilde{\alpha}_2}^{\text{ddNTK-I}} + m_{i;\tilde{\alpha}_2}^{\text{ddNTK-II}} \right) \end{aligned}$$

- ▶ Each network fits the training data with its own *particular* NTK, and so the resulting fully-trained particular output depends on the **NTK fluctuation**.
- ▶ The dependence on the **NTK differentials** means there's nontrivial representation learning at finite width.

Mean Prediction: *NTK Variance*

The fluctuation of the NTK term gives

$$m_{i;\delta}^{\Delta\text{NTK}} \equiv \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_4} \left(A_{(\delta\tilde{\alpha}_1)(\tilde{\alpha}_2\tilde{\alpha}_3)}^{(L)} + B_{\delta\tilde{\alpha}_2\tilde{\alpha}_1\tilde{\alpha}_3}^{(L)} + n_L B_{\delta\tilde{\alpha}_3\tilde{\alpha}_1\tilde{\alpha}_2}^{(L)} \right) \\ \times \tilde{H}_{(L)}^{\tilde{\alpha}_1\tilde{\alpha}_2} \tilde{H}_{(L)}^{\tilde{\alpha}_3\tilde{\alpha}_4} y_{i;\tilde{\alpha}_4},$$

where we decomposed the *NTK variance* into $A^{(L)}$ and $B^{(L)}$:

$$\mathbb{E} \left[\widehat{\Delta H}_{i_1 i_2; \alpha_1 \alpha_2}^{(L)} \widehat{\Delta H}_{i_3 i_4; \alpha_3 \alpha_4}^{(L)} \right] \\ \equiv \frac{1}{n_{L-1}} \left[\delta_{i_1 i_2} \delta_{i_3 i_4} A_{(\alpha_1 \alpha_2)(\alpha_3 \alpha_4)}^{(L)} + \delta_{i_1 i_3} \delta_{i_2 i_4} B_{\alpha_1 \alpha_3 \alpha_2 \alpha_4}^{(L)} + \delta_{i_1 i_4} \delta_{i_2 i_3} B_{\alpha_1 \alpha_4 \alpha_2 \alpha_3}^{(L)} \right].$$

Mean Prediction: *NTK Variance*

The fluctuation of the NTK term gives

$$m_{i;\delta}^{\Delta\text{NTK}} \equiv \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_4} \left(A_{(\delta\tilde{\alpha}_1)(\tilde{\alpha}_2\tilde{\alpha}_3)}^{(L)} + B_{\delta\tilde{\alpha}_2\tilde{\alpha}_1\tilde{\alpha}_3}^{(L)} + n_L B_{\delta\tilde{\alpha}_3\tilde{\alpha}_1\tilde{\alpha}_2}^{(L)} \right) \\ \times \tilde{H}_{(L)}^{\tilde{\alpha}_1\tilde{\alpha}_2} \tilde{H}_{(L)}^{\tilde{\alpha}_3\tilde{\alpha}_4} y_{i;\tilde{\alpha}_4},$$

where we decomposed the *NTK variance* into $A^{(L)}$ and $B^{(L)}$:

$$\mathbb{E} \left[\widehat{\Delta H}_{i_1 i_2; \alpha_1 \alpha_2}^{(L)} \widehat{\Delta H}_{i_3 i_4; \alpha_3 \alpha_4}^{(L)} \right] \\ \equiv \frac{1}{n_{L-1}} \left[\delta_{i_1 i_2} \delta_{i_3 i_4} A_{(\alpha_1 \alpha_2)(\alpha_3 \alpha_4)}^{(L)} + \delta_{i_1 i_3} \delta_{i_2 i_4} B_{\alpha_1 \alpha_3 \alpha_2 \alpha_4}^{(L)} + \delta_{i_1 i_4} \delta_{i_2 i_3} B_{\alpha_1 \alpha_4 \alpha_2 \alpha_3}^{(L)} \right].$$

Mean Prediction: *dNTK-Preactivation Cross Correlation*

The *dNTK* term gives

$$\begin{aligned}
 & m_{i;\delta}^{\text{dNTK}} \\
 \equiv & - \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_4} \left[2 \left(P_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} + Q_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} + n_L Q_{\delta \tilde{\alpha}_2 \tilde{\alpha}_1 \tilde{\alpha}_3}^{(L)} \right) Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \right. \\
 & \quad + \left(n_L P_{\tilde{\alpha}_1 \delta \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} + Q_{\tilde{\alpha}_1 \delta \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} + Q_{\tilde{\alpha}_1 \tilde{\alpha}_2 \delta \tilde{\alpha}_3}^{(L)} \right) Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \\
 & \quad \left. + \left(P_{\tilde{\alpha}_1 \delta \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} + n_L Q_{\tilde{\alpha}_1 \delta \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} + Q_{\tilde{\alpha}_1 \tilde{\alpha}_2 \delta \tilde{\alpha}_3}^{(L)} \right) Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_4 \tilde{\alpha}_3} \right] y_{i;\tilde{\alpha}_4},
 \end{aligned}$$

where we decomposed the *dNTK-preactivation cross correlators* into tensors $P^{(L)}$ and $Q^{(L)}$:

$$\begin{aligned}
 & \mathbb{E} \left[\widehat{dH}_{i_0 i_1 i_2; \delta_0 \delta_1 \delta_2}^{(L)} z_{i_3; \delta_3}^{(L)} \right] \\
 \equiv & \frac{1}{n_{L-1}} \left[\delta_{i_0 i_3} \delta_{i_1 i_2} P_{\delta_0 \delta_1 \delta_2 \delta_3}^{(L)} + \delta_{i_0 i_1} \delta_{i_2 i_3} Q_{\delta_0 \delta_1 \delta_2 \delta_3}^{(L)} + \delta_{i_0 i_2} \delta_{i_1 i_3} Q_{\delta_0 \delta_2 \delta_1 \delta_3}^{(L)} \right].
 \end{aligned}$$

Mean Prediction: *dNTK-Preactivation Cross Correlation*

The *dNTK* term gives

$$\begin{aligned}
 & m_{i;\delta}^{\text{dNTK}} \\
 \equiv & - \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_4} \left[2 \left(P_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} + Q_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} + n_L Q_{\delta \tilde{\alpha}_2 \tilde{\alpha}_1 \tilde{\alpha}_3}^{(L)} \right) Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \right. \\
 & \quad + \left(n_L P_{\tilde{\alpha}_1 \delta \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} + Q_{\tilde{\alpha}_1 \delta \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} + Q_{\tilde{\alpha}_1 \tilde{\alpha}_2 \delta \tilde{\alpha}_3}^{(L)} \right) Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \\
 & \quad \left. + \left(P_{\tilde{\alpha}_1 \delta \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} + n_L Q_{\tilde{\alpha}_1 \delta \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} + Q_{\tilde{\alpha}_1 \tilde{\alpha}_2 \delta \tilde{\alpha}_3}^{(L)} \right) Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_4 \tilde{\alpha}_3} \right] y_{i;\tilde{\alpha}_4},
 \end{aligned}$$

where we decomposed the *dNTK-preactivation cross correlators* into tensors $P^{(L)}$ and $Q^{(L)}$:

$$\begin{aligned}
 & \mathbb{E} \left[\widehat{dH}_{i_0 i_1 i_2; \delta_0 \delta_1 \delta_2}^{(L)} z_{i_3; \delta_3}^{(L)} \right] \\
 \equiv & \frac{1}{n_{L-1}} \left[\delta_{i_0 i_3} \delta_{i_1 i_2} P_{\delta_0 \delta_1 \delta_2 \delta_3}^{(L)} + \delta_{i_0 i_1} \delta_{i_2 i_3} Q_{\delta_0 \delta_1 \delta_2 \delta_3}^{(L)} + \delta_{i_0 i_2} \delta_{i_1 i_3} Q_{\delta_0 \delta_2 \delta_1 \delta_3}^{(L)} \right].
 \end{aligned}$$

Mean Prediction: $ddNTK_I$ Mean

The first $ddNTK$ term gives

$$\begin{aligned}
 m_{i;\delta}^{\text{ddNTK-I}} &\equiv \\
 &- \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_6} \left[R_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} \left(Z_{\text{IB}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} + Z_{\text{IB}}^{\tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_1 \tilde{\alpha}_5 \tilde{\alpha}_6 \tilde{\alpha}_4} + Z_{\text{IB}}^{\tilde{\alpha}_3 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_6 \tilde{\alpha}_4 \tilde{\alpha}_5} \right) + \right. \\
 &R_{\tilde{\alpha}_1 \delta \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} Z_{\text{IA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} + R_{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \delta}^{(L)} Z_{\text{IA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_5 \tilde{\alpha}_6 \tilde{\alpha}_4} + R_{\tilde{\alpha}_1 \tilde{\alpha}_3 \delta \tilde{\alpha}_2}^{(L)} Z_{\text{IA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_6 \tilde{\alpha}_4 \tilde{\alpha}_5} \\
 &\left. \right] \times \left[y_{i;\tilde{\alpha}_4} \left(\sum_j y_{j;\tilde{\alpha}_5} y_{j;\tilde{\alpha}_6} + n_L K_{\tilde{\alpha}_5 \tilde{\alpha}_6}^{(L)} \right) + y_{i;\tilde{\alpha}_5} K_{\tilde{\alpha}_6 \tilde{\alpha}_4}^{(L)} + y_{i;\tilde{\alpha}_6} K_{\tilde{\alpha}_4 \tilde{\alpha}_5}^{(L)} \right]
 \end{aligned}$$

where we decomposed the $dNTK_I$ mean into the tensor $R^{(L)}$:

$$\begin{aligned}
 &\mathbb{E} \left[\widehat{\text{dd}}_I H_{i_0 i_1 i_2 i_3; \delta_0 \delta_1 \delta_2 \delta_3}^{(L)} \right] \\
 &\equiv \frac{1}{n_{L-1}} \left[\delta_{i_0 i_1} \delta_{i_2 i_3} R_{\delta_0 \delta_1 \delta_2 \delta_3}^{(L)} + \delta_{i_0 i_2} \delta_{i_3 i_1} R_{\delta_0 \delta_2 \delta_3 \delta_1}^{(L)} + \delta_{i_0 i_3} \delta_{i_1 i_2} R_{\delta_0 \delta_3 \delta_1 \delta_2}^{(L)} \right].
 \end{aligned}$$

Mean Prediction: $ddNTK_I$ Mean

The first $ddNTK$ term gives

$$\begin{aligned}
 m_{i;\delta}^{\text{ddNTK-I}} &\equiv \\
 &- \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_6} \left[R_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} \left(Z_{\text{IB}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} + Z_{\text{IB}}^{\tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_1 \tilde{\alpha}_5 \tilde{\alpha}_6 \tilde{\alpha}_4} + Z_{\text{IB}}^{\tilde{\alpha}_3 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_6 \tilde{\alpha}_4 \tilde{\alpha}_5} \right) + \right. \\
 &R_{\tilde{\alpha}_1 \delta \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} Z_{\text{IA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} + R_{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \delta}^{(L)} Z_{\text{IA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_5 \tilde{\alpha}_6 \tilde{\alpha}_4} + R_{\tilde{\alpha}_1 \tilde{\alpha}_3 \delta \tilde{\alpha}_2}^{(L)} Z_{\text{IA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_6 \tilde{\alpha}_4 \tilde{\alpha}_5} \\
 &\left. \right] \times \left[y_{i;\tilde{\alpha}_4} \left(\sum_j y_{j;\tilde{\alpha}_5} y_{j;\tilde{\alpha}_6} + n_L K_{\tilde{\alpha}_5 \tilde{\alpha}_6}^{(L)} \right) + y_{i;\tilde{\alpha}_5} K_{\tilde{\alpha}_6 \tilde{\alpha}_4}^{(L)} + y_{i;\tilde{\alpha}_6} K_{\tilde{\alpha}_4 \tilde{\alpha}_5}^{(L)} \right]
 \end{aligned}$$

where we decomposed the $dNTK_I$ mean into the tensor $R^{(L)}$:

$$\begin{aligned}
 &\mathbb{E} \left[\widehat{\text{dd}}_I H_{i_0 i_1 i_2 i_3; \delta_0 \delta_1 \delta_2 \delta_3}^{(L)} \right] \\
 &\equiv \frac{1}{n_{L-1}} \left[\delta_{i_0 i_1} \delta_{i_2 i_3} R_{\delta_0 \delta_1 \delta_2 \delta_3}^{(L)} + \delta_{i_0 i_2} \delta_{i_3 i_1} R_{\delta_0 \delta_2 \delta_3 \delta_1}^{(L)} + \delta_{i_0 i_3} \delta_{i_1 i_2} R_{\delta_0 \delta_3 \delta_1 \delta_2}^{(L)} \right].
 \end{aligned}$$

Mean Prediction: $ddNNTK_{||}$ Mean

The second ddNTK term gives

$$\begin{aligned}
 & m_{i;\delta}^{\text{ddNTK-II}} \\
 \equiv & - \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_6} \left[S_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} Z_{\text{IIB}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} + T_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} Z_{\text{IIB}}^{\tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_1 \tilde{\alpha}_5 \tilde{\alpha}_6 \tilde{\alpha}_4} \right. \\
 & \quad + U_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} Z_{\text{IIB}}^{\tilde{\alpha}_3 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_6 \tilde{\alpha}_4 \tilde{\alpha}_5} + S_{\tilde{\alpha}_1 \tilde{\alpha}_2 \delta \tilde{\alpha}_3}^{(L)} Z_{\text{IIA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_5 \tilde{\alpha}_6 \tilde{\alpha}_4} \\
 & \quad \left. + T_{\tilde{\alpha}_1 \delta \tilde{\alpha}_3 \tilde{\alpha}_2}^{(L)} Z_{\text{IIA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} + U_{\tilde{\alpha}_1 \tilde{\alpha}_3 \tilde{\alpha}_2 \delta}^{(L)} Z_{\text{IIA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_6 \tilde{\alpha}_4 \tilde{\alpha}_5} \right] \\
 & \times \left[y_{i;\tilde{\alpha}_4} \left(\sum_j y_{j;\tilde{\alpha}_5} y_{j;\tilde{\alpha}_6} + n_L K_{\tilde{\alpha}_5 \tilde{\alpha}_6}^{(L)} \right) + y_{i;\tilde{\alpha}_5} K_{\tilde{\alpha}_6 \tilde{\alpha}_4}^{(L)} + y_{i;\tilde{\alpha}_6} K_{\tilde{\alpha}_4 \tilde{\alpha}_5}^{(L)} \right]
 \end{aligned}$$

where we decomposed the $dNNTK_{||}$ mean into $S^{(L)}$, $T^{(L)}$, $U^{(L)}$:

$$\begin{aligned}
 & \mathbb{E} \left[\widehat{dd_{||}H}_{i_1 i_2 i_3 i_4; \delta_1 \delta_2 \delta_3 \delta_4}^{(L)} \right] \\
 \equiv & \frac{1}{n_{L-1}} \left[\delta_{i_1 i_2} \delta_{i_3 i_4} S_{\delta_1 \delta_2 \delta_3 \delta_4}^{(L)} + \delta_{i_1 i_3} \delta_{i_4 i_2} T_{\delta_1 \delta_3 \delta_4 \delta_2}^{(L)} + \delta_{i_1 i_4} \delta_{i_2 i_3} U_{\delta_1 \delta_4 \delta_2 \delta_3}^{(L)} \right].
 \end{aligned}$$

Mean Prediction: $ddNTK_{||}$ Mean

The second ddNTK term gives

$$\begin{aligned}
 & m_{i;\delta}^{\text{ddNTK-II}} \\
 \equiv & - \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_6} \left[S_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} Z_{\text{IIB}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} + T_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} Z_{\text{IIB}}^{\tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_1 \tilde{\alpha}_5 \tilde{\alpha}_6 \tilde{\alpha}_4} \right. \\
 & \quad + U_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} Z_{\text{IIB}}^{\tilde{\alpha}_3 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_6 \tilde{\alpha}_4 \tilde{\alpha}_5} + S_{\tilde{\alpha}_1 \tilde{\alpha}_2 \delta \tilde{\alpha}_3}^{(L)} Z_{\text{IIA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_5 \tilde{\alpha}_6 \tilde{\alpha}_4} \\
 & \quad \left. + T_{\tilde{\alpha}_1 \delta \tilde{\alpha}_3 \tilde{\alpha}_2}^{(L)} Z_{\text{IIA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} + U_{\tilde{\alpha}_1 \tilde{\alpha}_3 \tilde{\alpha}_2 \delta}^{(L)} Z_{\text{IIA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_6 \tilde{\alpha}_4 \tilde{\alpha}_5} \right] \\
 & \times \left[y_{i;\tilde{\alpha}_4} \left(\sum_j y_{j;\tilde{\alpha}_5} y_{j;\tilde{\alpha}_6} + n_L K_{\tilde{\alpha}_5 \tilde{\alpha}_6}^{(L)} \right) + y_{i;\tilde{\alpha}_5} K_{\tilde{\alpha}_6 \tilde{\alpha}_4}^{(L)} + y_{i;\tilde{\alpha}_6} K_{\tilde{\alpha}_4 \tilde{\alpha}_5}^{(L)} \right]
 \end{aligned}$$

where we decomposed the $dNTK_{||}$ mean into $S^{(L)}$, $T^{(L)}$, $U^{(L)}$:

$$\begin{aligned}
 & \mathbb{E} \left[\widehat{dd_{||}H}_{i_1 i_2 i_3 i_4; \delta_1 \delta_2 \delta_3 \delta_4}^{(L)} \right] \\
 \equiv & \frac{1}{n_{L-1}} \left[\delta_{i_1 i_2} \delta_{i_3 i_4} S_{\delta_1 \delta_2 \delta_3 \delta_4}^{(L)} + \delta_{i_1 i_3} \delta_{i_4 i_2} T_{\delta_1 \delta_3 \delta_4 \delta_2}^{(L)} + \delta_{i_1 i_4} \delta_{i_2 i_3} U_{\delta_1 \delta_4 \delta_2 \delta_3}^{(L)} \right].
 \end{aligned}$$

Prediction Variance

In addition to the ensemble mean, we can consider other statistics:

$$\begin{aligned} & \text{Cov} \left[z_{i_1; \hat{\beta}_1}^{(L)}(\infty), z_{i_2; \hat{\beta}_2}^{(L)}(\infty) \right] \\ & \equiv \mathbb{E} \left[z_{i_1; \hat{\beta}_1}^{(L)}(\infty) z_{i_2; \hat{\beta}_2}^{(L)}(\infty) \right] - \mathbb{E} \left[z_{i_1; \hat{\beta}_1}^{(L)}(\infty) \right] \mathbb{E} \left[z_{i_2; \hat{\beta}_2}^{(L)}(\infty) \right]. \end{aligned}$$

Prediction Variance

In addition to the ensemble mean, we can consider other statistics:

$$\begin{aligned} & \text{Cov} \left[z_{i_1; \hat{\beta}_1}^{(L)}(\infty), z_{i_2; \hat{\beta}_2}^{(L)}(\infty) \right] \\ & \equiv \mathbb{E} \left[z_{i_1; \hat{\beta}_1}^{(L)}(\infty) z_{i_2; \hat{\beta}_2}^{(L)}(\infty) \right] - \mathbb{E} \left[z_{i_1; \hat{\beta}_1}^{(L)}(\infty) \right] \mathbb{E} \left[z_{i_2; \hat{\beta}_2}^{(L)}(\infty) \right]. \end{aligned}$$

While we won't print this quantity in full – the full expression doesn't really play nicely with the constraints of the slides – you can easily extract insight by considering specific contributions.

Prediction Variance

To see another manifestation of output “wiring,” consider

$$\sum_{\substack{j_1, j_2 \\ \tilde{\alpha}_1, \dots, \tilde{\alpha}_4}} \mathbb{E} \left[z_{i_2; \hat{\beta}_2}^{(L)} \widehat{dH}_{i_1 j_1 j_2; \hat{\beta}_1 \tilde{\alpha}_1 \tilde{\alpha}_2}^{(L)} \right] z_{\mathbf{B}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \\ \times \mathbb{E} \left[\left(z_{j_1; \tilde{\alpha}_3}^{(L)} - y_{j_1; \tilde{\alpha}_3} \right) \left(z_{j_2; \tilde{\alpha}_4}^{(L)} - y_{j_2; \tilde{\alpha}_4} \right) \right]$$

Prediction Variance

To see another manifestation of output “wiring,” consider

$$\begin{aligned}
 & \sum_{\substack{j_1, j_2 \\ \tilde{\alpha}_1, \dots, \tilde{\alpha}_4}} \mathbb{E} \left[z_{i_2; \dot{\beta}_2}^{(L)} \widehat{dH}_{i_1 j_1 j_2; \dot{\beta}_1 \tilde{\alpha}_1 \tilde{\alpha}_2}^{(L)} \right] Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \\
 & \quad \times \mathbb{E} \left[\left(z_{j_1; \tilde{\alpha}_3}^{(L)} - y_{j_1; \tilde{\alpha}_3} \right) \left(z_{j_2; \tilde{\alpha}_4}^{(L)} - y_{j_2; \tilde{\alpha}_4} \right) \right] \\
 &= \frac{1}{n_{L-1}} \delta_{i_1 i_2} \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_4} \left[P_{\dot{\beta}_1 \tilde{\alpha}_1 \tilde{\alpha}_2 \dot{\beta}_2}^{(L)} \left(\sum_j y_{j; \tilde{\alpha}_3} y_{j; \tilde{\alpha}_4} \right) \right. \\
 & \quad \left. + \left(n_L P_{\dot{\beta}_1 \tilde{\alpha}_1 \tilde{\alpha}_2 \dot{\beta}_2}^{(L)} + Q_{\dot{\beta}_1 \tilde{\alpha}_1 \tilde{\alpha}_2 \dot{\beta}_2}^{(L)} + Q_{\dot{\beta}_1 \tilde{\alpha}_2 \tilde{\alpha}_1 \dot{\beta}_2}^{(L)} \right) G_{\tilde{\alpha}_3 \tilde{\alpha}_4}^{(L)} \right] Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \\
 & \quad + \frac{1}{n_{L-1}} \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_4} Q_{\dot{\beta}_1 \tilde{\alpha}_1 \tilde{\alpha}_2 \dot{\beta}_2}^{(L)} Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \left(y_{i_1; \tilde{\alpha}_3} y_{i_2; \tilde{\alpha}_4} + y_{i_1; \tilde{\alpha}_4} y_{i_2; \tilde{\alpha}_3} \right),
 \end{aligned}$$

Prediction Variance

To see another manifestation of output “wiring,” consider

$$\begin{aligned}
 & \sum_{\substack{j_1, j_2 \\ \tilde{\alpha}_1, \dots, \tilde{\alpha}_4}} \mathbb{E} \left[z_{i_2; \dot{\beta}_2}^{(L)} \widehat{dH}_{i_1 j_1 j_2; \dot{\beta}_1 \tilde{\alpha}_1 \tilde{\alpha}_2}^{(L)} \right] Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \\
 & \quad \times \mathbb{E} \left[\left(z_{j_1; \tilde{\alpha}_3}^{(L)} - y_{j_1; \tilde{\alpha}_3} \right) \left(z_{j_2; \tilde{\alpha}_4}^{(L)} - y_{j_2; \tilde{\alpha}_4} \right) \right] \\
 &= \frac{1}{n_{L-1}} \delta_{i_1 i_2} \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_4} \left[P_{\dot{\beta}_1 \tilde{\alpha}_1 \tilde{\alpha}_2 \dot{\beta}_2}^{(L)} \left(\sum_j y_{j; \tilde{\alpha}_3} y_{j; \tilde{\alpha}_4} \right) \right. \\
 & \quad \left. + \left(n_L P_{\dot{\beta}_1 \tilde{\alpha}_1 \tilde{\alpha}_2 \dot{\beta}_2}^{(L)} + Q_{\dot{\beta}_1 \tilde{\alpha}_1 \tilde{\alpha}_2 \dot{\beta}_2}^{(L)} + Q_{\dot{\beta}_1 \tilde{\alpha}_2 \tilde{\alpha}_1 \dot{\beta}_2}^{(L)} \right) G_{\tilde{\alpha}_3 \tilde{\alpha}_4}^{(L)} \right] Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \\
 & \quad + \frac{1}{n_{L-1}} \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_4} Q_{\dot{\beta}_1 \tilde{\alpha}_1 \tilde{\alpha}_2 \dot{\beta}_2}^{(L)} Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \left(y_{i_1; \tilde{\alpha}_3} y_{i_2; \tilde{\alpha}_4} + y_{i_1; \tilde{\alpha}_4} y_{i_2; \tilde{\alpha}_3} \right),
 \end{aligned}$$

- *Wiring* is exhibited when $y_{i_1} \neq 0$ and $y_{i_2} \neq 0$.

Prediction is Nearly-Gaussian

At finite width, $p(z^{(L)}(\infty) | \mathcal{D})$ now has non-Gaussian statistics:

Prediction is Nearly-Gaussian

At finite width, $p(z^{(L)}(\infty) | \mathcal{D})$ now has non-Gaussian statistics:

$$\blacktriangleright z_{i;\hat{\beta}}(\infty) \equiv z_{i;\hat{\beta}}(\infty) \left[z^{(L)}, \widehat{H}^{(L)}, \widehat{dH}^{(L)}, \widehat{dd}_I H^{(L)}, \widehat{dd}_{II} H^{(L)} \right].$$

Prediction is Nearly-Gaussian

At finite width, $p(z^{(L)}(\infty) | \mathcal{D})$ now has non-Gaussian statistics:

- ▶ $z_{i;\hat{\beta}}(\infty) \equiv z_{i;\hat{\beta}}(\infty) \left[z^{(L)}, \hat{H}^{(L)}, \widehat{dH}^{(L)}, \widehat{dd}_I H^{(L)}, \widehat{dd}_{II} H^{(L)} \right]$.
- ▶ $p(z^{(L)}, \hat{H}^{(L)}, \widehat{dH}^{(L)}, \widehat{dd}_I H^{(L)}, \widehat{dd}_{II} H^{(L)} | \mathcal{D})$ is *nearly-Gaussian*.

Prediction is Nearly-Gaussian

At finite width, $p(z^{(L)}(\infty) | \mathcal{D})$ now has non-Gaussian statistics:

- ▶ $z_{i;\hat{\beta}}(\infty) \equiv z_{i;\hat{\beta}}(\infty) \left[z^{(L)}, \widehat{H}^{(L)}, \widehat{dH}^{(L)}, \widehat{dd}_I H^{(L)}, \widehat{dd}_{II} H^{(L)} \right]$.
- ▶ $p(z^{(L)}, \widehat{H}^{(L)}, \widehat{dH}^{(L)}, \widehat{dd}_I H^{(L)}, \widehat{dd}_{II} H^{(L)} | \mathcal{D})$ is *nearly-Gaussian*.
- ▶ Explicit expressions of higher-point correlators are challenging to display in any media format...

Generalization

The **generalization error** is a quantitative measure of how well a network is really approximating the desired function:

$$\mathcal{E} \equiv \mathcal{L}_{\mathcal{B}} - \mathcal{L}_{\mathcal{A}}.$$

Generalization

The **generalization error** is a quantitative measure of how well a network is really approximating the desired function:

$$\mathcal{E} \equiv \mathcal{L}_{\mathcal{B}}.$$

Generalization

The **generalization error** is a quantitative measure of how well a network is really approximating the desired function:

$$\mathcal{E} \equiv \mathcal{L}_{\mathcal{B}}.$$

Let's evaluate the **MSE test loss**, averaged over an ensemble of fully-trained networks:

$$\mathbb{E} [\mathcal{L}_{\mathcal{B}}(T)] = \mathbb{E} \left[\frac{1}{2} \sum_{i=1}^{n_L} \sum_{\hat{\beta} \in \mathcal{B}} \left(z_{i;\hat{\beta}}^{(L)}(\infty) - y_{i;\hat{\beta}} \right)^2 \right].$$

Generalization

The **generalization error** is a quantitative measure of how well a network is really approximating the desired function:

$$\mathcal{E} \equiv \mathcal{L}_{\mathcal{B}}.$$

Let's evaluate the **MSE test loss**, averaged over an ensemble of fully-trained networks:

$$\mathbb{E}[\mathcal{L}_{\mathcal{B}}(T)] = \frac{1}{2} \sum_{\dot{\beta} \in \mathcal{B}} \left\{ \sum_{i=1}^{n_L} (m_{i;\dot{\beta}} - y_{i;\dot{\beta}})^2 + \sum_{i=1}^{n_L} \text{Cov}[z_{i;\dot{\beta}}^{(L)}(\infty), z_{i;\dot{\beta}}^{(L)}(\infty)] \right\}.$$

Generalization

The **generalization error** is a quantitative measure of how well a network is really approximating the desired function:

$$\mathcal{E} \equiv \mathcal{L}_{\mathcal{B}}.$$

Let's evaluate the **MSE test loss**, averaged over an ensemble of fully-trained networks:

$$\mathbb{E}[\mathcal{L}_{\mathcal{B}}(T)] = \frac{1}{2} \sum_{\dot{\beta} \in \mathcal{B}} \left\{ \sum_{i=1}^{n_L} (m_{i;\dot{\beta}} - y_{i;\dot{\beta}})^2 + \sum_{i=1}^{n_L} \text{Cov}[z_{i;\dot{\beta}}^{(L)}(\infty), z_{i;\dot{\beta}}^{(L)}(\infty)] \right\}.$$

- ▶ This illustrates a **generalized bias-variance tradeoff**: different settings of the hyperparameters will decrease one term at the cost of increasing the other.

Generalization

The **generalization error** is a quantitative measure of how well a network is really approximating the desired function:

$$\mathcal{E} \equiv \mathcal{L}_{\mathcal{B}}.$$

Let's evaluate the **MSE test loss**, averaged over an ensemble of fully-trained networks:

$$\mathbb{E}[\mathcal{L}_{\mathcal{B}}(T)] = \frac{1}{2} \sum_{\hat{\beta} \in \mathcal{B}} \left\{ \sum_{i=1}^{n_L} (m_{i;\hat{\beta}} - y_{i;\hat{\beta}})^2 + \sum_{i=1}^{n_L} \text{Cov}[z_{i;\hat{\beta}}^{(L)}(\infty), z_{i;\hat{\beta}}^{(L)}(\infty)] \right\}.$$

- ▶ This illustrates a **generalized bias-variance tradeoff**: different settings of the hyperparameters will decrease one term at the cost of increasing the other.
- ▶ Importantly, we're choosing between *ensembles* not *models*.

Generalization: Interpolation and Extrapolation

Given the true outputs $y_{i;\pm}$ for *two* inputs $x_{i;\pm} = x_{i;0} \pm \frac{\delta x_i}{2}$, what is the prediction for a one-parameter family of test inputs,

$$s x_{i;+} + (1 - s) x_{i;-} = x_{i;0} + \frac{(2s - 1)}{2} \delta x_i \equiv x_{i;(2s-1)},$$

that sit on a line passing through $x_{i;+}$ and $x_{i;-}$?

Generalization: Interpolation and Extrapolation

Given the true outputs $y_{i;\pm}$ for *two* inputs $x_{i;\pm} = x_{i;0} \pm \frac{\delta x_i}{2}$, what is the prediction for a one-parameter family of test inputs,

$$s x_{i;+} + (1 - s) x_{i;-} = x_{i;0} + \frac{(2s - 1)}{2} \delta x_i \equiv x_{i;(2s-1)},$$

that sit on a line passing through $x_{i;+}$ and $x_{i;-}$?

- ▶ When our parameter s is inside the unit interval $s \in [0, 1]$, this is a question about neural-network **interpolation**.
- ▶ For s outside the unit interval, it's **extrapolation**.
- ▶ For general s , let's refer to this collectively as ***-polation**.

*-polation at Infinite Width

To understand *-**polation** at infinite by *smooth activations*, we need to evaluate:

$$z_{i;\hat{\beta}}^{(L)}(\infty) = z_{i;\hat{\beta}}^{(L)} - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2 \in \mathcal{A}} \Theta_{\hat{\beta}\tilde{\alpha}_1}^{(L)} \tilde{\Theta}_{(L)}^{\tilde{\alpha}_1\tilde{\alpha}_2} \left(z_{i;\tilde{\alpha}_2}^{(L)} - y_{i;\tilde{\alpha}_2} \right) .$$

*-polation at Infinite Width

To understand *-**polation** at infinite by *smooth activations*, we need to evaluate:

$$z_{i;\hat{\beta}}^{(L)}(\infty) = z_{i;\hat{\beta}}^{(L)} - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2 \in \mathcal{A}} \Theta_{\hat{\beta}\tilde{\alpha}_1}^{(L)} \tilde{\Theta}_{(\tilde{\alpha}_1\tilde{\alpha}_2)}^{(L)} \left(z_{i;\tilde{\alpha}_2}^{(L)} - y_{i;\tilde{\alpha}_2} \right).$$

- ▶ We need to invert the two-by-two submatrix of the NTK on the training set only, $\tilde{\Theta}_{(\tilde{\alpha}_1\tilde{\alpha}_2)}^{(L)}$.
- ▶ We will need to evaluate elements of the NTK between our test and training set, $\Theta_{(2s-1)\pm}^{(L)}$.

*-polation at Infinite Width

To understand *-**polation** at infinite by *smooth activations*, we need to evaluate:

$$z_{i;\hat{\beta}}^{(L)}(\infty) = z_{i;\hat{\beta}}^{(L)} - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2 \in \mathcal{A}} \Theta_{\hat{\beta}\tilde{\alpha}_1}^{(L)} \tilde{\Theta}_{(\tilde{\alpha}_1\tilde{\alpha}_2)}^{(L)} \left(z_{i;\tilde{\alpha}_2}^{(L)} - y_{i;\tilde{\alpha}_2} \right).$$

The inverse of the submatrix can be written as:

$$\tilde{\Theta}^{\tilde{\alpha}_1\tilde{\alpha}_2} = \frac{1}{\Theta_{++}\Theta_{--} - \Theta_{+-}^2} \begin{pmatrix} \Theta_{--} & -\Theta_{+-} \\ -\Theta_{+-} & \Theta_{++} \end{pmatrix}.$$

*-polation at Infinite Width

To understand *-**polation** at infinite by *smooth activations*, we need to evaluate:

$$z_{i;\dot{\beta}}^{(L)}(\infty) = z_{i;\dot{\beta}}^{(L)} - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2 \in \mathcal{A}} \Theta_{\dot{\beta}\tilde{\alpha}_1}^{(L)} \tilde{\Theta}_{(L)}^{\tilde{\alpha}_1\tilde{\alpha}_2} \left(z_{i;\tilde{\alpha}_2}^{(L)} - y_{i;\tilde{\alpha}_2} \right).$$

The test-train NTK can be evaluated as

$$\Theta_{(2s-1)\pm} = s\Theta_{\pm\pm} + (1-s)\Theta_{\pm-} - 2s(1-s)\delta\delta\Theta_{[0]} + O(\delta^3),$$

where $\delta\delta\Theta_{[0]}$ is the difference

$$\delta\delta\Theta_{[0]} \equiv \frac{1}{4} \left[\Theta_{++} + \Theta_{--} + 2\Theta_{+-} \right] - \Theta_{00} + O(\delta^4).$$

*-polation at Infinite Width

As an illustration, consider two training inputs with the same norm:

$$\begin{aligned} z_{i;(2s-1)}^{(L)}(\infty) &= \left[z_{i;(2s-1)}^{(L)} - s z_{i;+}^{(L)} - (1-s) z_{i;-}^{(L)} \right] + [s y_{i;+} + (1-s) y_{i;-}] \\ &\quad - 4s(1-s) \left(\frac{\delta \delta \Theta_{[0]}}{\Theta_{00}} \right) (z_{i;+}^{(L)} + z_{i;-}^{(L)} + y_{i;+} + y_{i;-}) \\ &\quad + O(\delta^3) . \end{aligned}$$

*-polation at Infinite Width

As an illustration, consider two training inputs with the same norm:

$$\begin{aligned} z_{i;(2s-1)}^{(L)}(\infty) &= \left[z_{i;(2s-1)}^{(L)} - s z_{i;+}^{(L)} - (1-s) z_{i;-}^{(L)} \right] + [s y_{i;+} + (1-s) y_{i;-}] \\ &\quad - 4s(1-s) \left(\frac{\delta \delta \Theta_{[0]}}{\Theta_{00}} \right) (z_{i;+}^{(L)} + z_{i;-}^{(L)} + y_{i;+} + y_{i;-}) \\ &\quad + O(\delta^3) . \end{aligned}$$

- ▶ *Nonlinear* networks can *nonlinearly* *-polate!

*-polation at Infinite Width

The prediction of the ensemble has a **bias**:

$$m_{i;(2s-1)}^{\infty} - y_{i;(2s-1)} = \left[y_{i;+} + (1-s)y_{i;-} - y_{i;(2s-1)} \right] \\ - 4s(1-s) \left(\frac{\delta \delta \Theta_{[0]}}{\Theta_{00}} \right) (y_{i;+} + y_{i;-}) + O(\delta^3) .$$

*-polation at Infinite Width

The prediction of the ensemble has a **bias**:

$$m_{i;(2s-1)}^{\infty} - y_{i;(2s-1)} = \left[y_{i;+} + (1-s)y_{i;-} - y_{i;(2s-1)} \right] \\ - 4s(1-s) \left(\frac{\delta \delta \Theta_{[0]}}{\Theta_{00}} \right) (y_{i;+} + y_{i;-}) + O(\delta^3).$$

- ▶ This decomposes into a part measuring the **nonlinearity in the labels** and the **network curvature** around $(y_{i;+} + y_{i;-})/2$.

*-polation at Infinite Width

The prediction of the ensemble has a **bias**:

$$m_{i;(2s-1)}^{\infty} - y_{i;(2s-1)} = \left[y_{i;+} + (1-s)y_{i;-} - y_{i;(2s-1)} \right] \\ - 4s(1-s) \left(\frac{\delta\delta\Theta_{[0]}}{\Theta_{00}} \right) (y_{i;+} + y_{i;-}) + O(\delta^3) .$$

- ▶ This decomposes into a part measuring the nonlinearity in the labels and the **network curvature** around $(y_{i;+} + y_{i;-})/2$.
- ▶ This curvature encodes the **inductive bias** of the function computed by the network.

*-polation at Finite Width

The prediction of the ensemble has a **bias**:

$$m_{i;(2s-1)} - y_{i;(2s-1)} = \left[y_{i;+} + (1-s)y_{i;-} - y_{i;(2s-1)} \right] + O(y^3) \\ - 4s(1-s) \left(\frac{\delta\delta\Theta_{[0]}}{\Theta_{00}} \right) (y_{i;+} + y_{i;-}) + O(\delta^3) .$$

- ▶ This decomposes into a part measuring the nonlinearity in the labels and the **network curvature** around $(y_{i;+} + y_{i;-})/2$.
- ▶ This curvature encodes the **inductive bias** of the function computed by the network.
- ▶ At *finite width*, we'd have a **cubic *-polation**.

Generalization at Finite Width

All terms are proportional to one of these dimensionless ratios:

$$\begin{array}{cccc} \frac{A^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^2}, & \frac{B^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^2}, & \frac{P^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^2}, & \frac{Q^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^2}, \\ \frac{R^{(L)} K^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^3}, & \frac{S^{(L)} K^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^3}, & \frac{T^{(L)} K^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^3}, & \frac{U^{(L)} K^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^3}. \end{array}$$

Generalization at Finite Width

All terms are proportional to one of these dimensionless ratios:

$$\begin{aligned} & \frac{A^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^2}, & \frac{B^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^2}, & \frac{P^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^2}, & \frac{Q^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^2}, \\ & \frac{R^{(L)} K^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^3}, & \frac{S^{(L)} K^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^3}, & \frac{T^{(L)} K^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^3}, & \frac{U^{(L)} K^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^3}. \end{aligned}$$

Overall we should find for the finite-width corrections:

$$m_{i;\dot{\beta}} - m_{i;\dot{\beta}}^{\infty} = O\left(\frac{L}{n}\right).$$

Optimal Aspect Ratio

In deep networks, we want to balance the **positive effect of *representation learning*** against the **negative effect of *fluctuations***.

Optimal Aspect Ratio

In deep networks, we want to balance the **positive effect of *representation learning*** against the **negative effect of *fluctuations***.

- ▶ We can look at the finite-width generalization error – and someone should(!) – but requires studying $O(L^2/n^2)$ effects.

Optimal Aspect Ratio

In deep networks, we want to balance the **positive effect of *representation learning*** against the **negative effect of *fluctuations***.

- ▶ We can look at the finite-width generalization error – and someone should(!) – but requires studying $O(L^2/n^2)$ effects.
- ▶ Can try to pick a simpler quantity – perhaps one defined at initialization – where we can compute higher-order effects.

Aside: Information Theory

The **entropy** of a probability distribution is given by

$$\mathcal{S}[p(x)] \equiv - \sum_x p(x) \log p(x).$$

Aside: Information Theory

The **entropy** of a probability distribution is given by

$$\mathcal{S}[p(x)] \equiv - \sum_x p(x) \log p(x).$$

- ▶ Entropy is *functional* of the distribution, taking a distribution as an argument and outputting a number.

Aside: Information Theory

The **entropy** of a probability distribution is given by

$$\mathcal{S}[p(x)] \equiv - \sum_x p(x) \log p(x).$$

- ▶ Entropy is *functional* of the distribution, taking a distribution as an argument and outputting a number.
- ▶ Quantitative measure of how much expected **information** is gained after making an *observation*.

Aside: Information Theory

The entropy is **additive** for two independent random variables x and y , with $p(x, y) = p(x)p(y)$:

$$\mathcal{S}[p(x, y)] = \mathcal{S}[p(x)] + \mathcal{S}[p(y)].$$

Aside: Information Theory

The entropy is **additive** for two independent random variables x and y , with $p(x, y) = p(x)p(y)$:

$$\begin{aligned} \mathcal{S}[p(x, y)] &= - \sum_{x, y} p(x, y) \log p(x, y) \\ &= - \sum_{x, y} p(x)p(y) [\log p(x) + \log p(y)] \\ &= - \sum_x p(x) \log p(x) - \sum_{y \in \mathcal{Y}} p(y) \log p(y) \\ &= \mathcal{S}[p(x)] + \mathcal{S}[p(y)] . \end{aligned}$$

Aside: Information Theory

The entropy is **additive** for two independent random variables x and y , with $p(x, y) = p(x)p(y)$:

$$\mathcal{S}[p(x, y)] = \mathcal{S}[p(x)] + \mathcal{S}[p(y)].$$

Aside: Information Theory

The entropy is **additive** for two independent random variables x and y , with $p(x, y) = p(x)p(y)$:

$$\mathcal{S}[p(x, y)] = \mathcal{S}[p(x)] + \mathcal{S}[p(y)] .$$

For two statistically *dependent* observables, constrained by a nonzero interaction, the entropy is **subadditive**:

$$\mathcal{S}[p(x, y)] < \mathcal{S}[p(x)] + \mathcal{S}[p(y)] .$$

Aside: Information Theory

The entropy is **additive** for two independent random variables x and y , with $p(x, y) = p(x)p(y)$:

$$\mathcal{S}[p(x, y)] = \mathcal{S}[p(x)] + \mathcal{S}[p(y)] .$$

For two statistically *dependent* observables, constrained by a nonzero interaction, the entropy is **subadditive**:

$$\mathcal{S}[p(x, y)] < \mathcal{S}[p(x)] + \mathcal{S}[p(y)] .$$

This is physically intuitive: for *dependent* variables observing y doesn't inform as much as it would have if we didn't know x .

Aside: Information Theory

The **mutual information** (MI) between two random variables is

$$\begin{aligned}\mathcal{I}[p(x, y)] &\equiv \mathcal{S}[p(x)] + \mathcal{S}[p(y)] - \mathcal{S}[p(x, y)] \\ &= \sum_{x, y} p(x, y) \log \left[\frac{p(x, y)}{p(x) p(y)} \right].\end{aligned}$$

Aside: Information Theory

The **mutual information** (MI) between two random variables is

$$\begin{aligned}\mathcal{I}[p(x, y)] &\equiv \mathcal{S}[p(x)] + \mathcal{S}[p(y)] - \mathcal{S}[p(x, y)] \\ &= \sum_{x, y} p(x, y) \log \left[\frac{p(x, y)}{p(x) p(y)} \right].\end{aligned}$$

- ▶ A functional of a joint probability distribution.

Aside: Information Theory

The **mutual information** (MI) between two random variables is

$$\begin{aligned}\mathcal{I}[p(x, y)] &\equiv \mathcal{S}[p(x)] + \mathcal{S}[p(y)] - \mathcal{S}[p(x, y)] \\ &= \sum_{x, y} p(x, y) \log \left[\frac{p(x, y)}{p(x) p(y)} \right].\end{aligned}$$

- ▶ A functional of a joint probability distribution.
- ▶ An average measure of how much information an observation of x conveys about an observation of y , and vice versa.

Aside: Information Theory

Rearranging, we see that the *subadditivity* of the entropy implies the nonnegativity of the mutual information:

$$\mathcal{I}[p(x, y)] \geq 0.$$

Aside: Information Theory

Rearranging, we see that the *subadditivity* of the entropy implies the nonnegativity of the mutual information:

$$\mathcal{I}[p(x, y)] \geq 0.$$

The MI of a joint distribution is telling us about the *interactions* that create the nontrivial non-Gaussian correlations, making it a *diagnostic of statistical dependence*.

Optimal Aspect Ratio

In deep networks, we want to balance the **positive effect of *representation learning*** against the **negative effect of *fluctuations***.

Optimal Aspect Ratio

In deep networks, we want to balance the **positive effect of *representation learning*** against the **negative effect of *fluctuations***.

- ▶ Intuitively, correlation between deep-layer neurons – or “wiring” – is a consequence of NTK differentials.

Optimal Aspect Ratio

In deep networks, we want to balance the **positive effect of representation learning** against the **negative effect of fluctuations**.

- ▶ Intuitively, correlation between deep-layer neurons – or “wiring” – is a consequence of NTK differentials.
- ▶ The MI could provide an accessible measure of the potential for such **inductive bias** in the *prior distribution*.

Optimal Aspect Ratio

Consider the *single-input* MI between two sets of neurons,
 $\mathcal{M}_1 = \{1, \dots, m_1\}$ and $\mathcal{M}_2 = \{m_1 + 1, \dots, m_1 + m_2\}$, in layer L :

$$\mathcal{I}[p(\mathcal{M}_1, \mathcal{M}_2|x)] \equiv \mathcal{S}[p(\mathcal{M}_1|x)] + \mathcal{S}[p(\mathcal{M}_2|x)] - \mathcal{S}[p(\mathcal{M}_1, \mathcal{M}_2|x)] .$$

Optimal Aspect Ratio

Consider the *single-input* MI between two sets of neurons, $\mathcal{M}_1 = \{1, \dots, m_1\}$ and $\mathcal{M}_2 = \{m_1 + 1, \dots, m_1 + m_2\}$, in layer L :

$$\mathcal{I}[p(\mathcal{M}_1, \mathcal{M}_2|x)] \equiv \mathcal{S}[p(\mathcal{M}_1|x)] + \mathcal{S}[p(\mathcal{M}_2|x)] - \mathcal{S}[p(\mathcal{M}_1, \mathcal{M}_2|x)] .$$

With $r \equiv L/n$, we can *estimate* the optimal aspect ratio as

$$r^* \equiv \arg \max_{\{r, \mathcal{M}_1, \mathcal{M}_2\}} \mathcal{I}[p(\mathcal{M}_1, \mathcal{M}_2|x)] .$$

Aside: Unsupervised Learning

Maximizing this MI is related to **unsupervised learning** objectives.

Aside: Unsupervised Learning

Maximizing this MI is related to **unsupervised learning** objectives.

- ▶ The **InfoMax principle** recommends maximizing the mutual information between the input x and a *representation* $z(x)$.
- ▶ A related notion involves maximizing the mutual information between different representations, $z_1(x)$ and $z_2(x)$, for the same input x . This latter notion can be shown to *lower bound* the InfoMax objective and thus motivates our analysis here.

Slide unavailable.

Third-Order Result: Entropy

Using a **variational principle**, we can compute \mathcal{S} to $O(1/n^3)$:

$$\begin{aligned} \mathcal{S}[p(z_1, \dots, z_m | x)] &= \frac{m}{2} \log(2\pi eG) - \frac{(m^2 + 2m)}{16} \left(\frac{V}{nG^2}\right)^2 \\ &\quad + \frac{(m^3 + 10m^2 + 16m)}{48} \left(\frac{V}{nG^2}\right)^3 + O\left(\frac{1}{n^4}\right). \end{aligned}$$

Third-Order Result: Entropy

Using a **variational principle**, we can compute \mathcal{S} to $O(1/n^3)$:

$$\begin{aligned}\mathcal{S}[p(z_1, \dots, z_m | x)] &= \frac{m}{2} \log(2\pi eG) - \frac{(m^2 + 2m)}{16} \left(\frac{V}{nG^2}\right)^2 \\ &\quad + \frac{(m^3 + 10m^2 + 16m)}{48} \left(\frac{V}{nG^2}\right)^3 + O\left(\frac{1}{n^4}\right).\end{aligned}$$

- Note that the correction is definitely negative.

Third-Order Result: Entropy

Using a **variational principle**, we can compute \mathcal{S} to $O(1/n^3)$:

$$\begin{aligned} \mathcal{S}[p(z_1, \dots, z_m|x)] &= \frac{m}{2} \log(2\pi eG) - \frac{(m^2 + 2m)}{16} \left(\frac{V}{nG^2}\right)^2 \\ &\quad + \frac{(m^3 + 10m^2 + 16m)}{48} \left(\frac{V}{nG^2}\right)^3 + O\left(\frac{1}{n^4}\right). \end{aligned}$$

- ▶ Note that the correction is definitely negative.
- ▶ Note that unlike all our previous results, the leading correction here is *second order* in the inverse layer width $\sim V^2/n^2$.

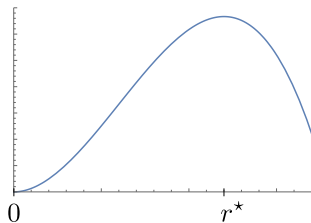
Optimal Aspect Ratio

Defining for the four-point vertex ratio,

$$\frac{V^{(L)}}{n_{L-1} (G^{(L)})^2} \equiv \nu r,$$

with ν an *activation function-dependent constant*, we get:

$$\mathcal{I}[p(\mathcal{M}_1, \mathcal{M}_2|x)] = \frac{m_1 m_2}{8} \nu^2 r^2 - \frac{m_1 m_2 (20 + 3m_1 + 3m_2)}{48} \nu^3 r^3 + O(r^4)$$



Optimal Aspect Ratio

Defining for the four-point vertex ratio,

$$\frac{V^{(L)}}{n_{L-1} (G^{(L)})^2} \equiv \nu r,$$

with ν an *activation function-dependent constant*, we get:

$$\mathcal{I} [p(\mathcal{M}_1, \mathcal{M}_2|x)] = \frac{m_1 m_2}{8} \nu^2 r^2 - \frac{m_1 m_2 (20 + 3m_1 + 3m_2)}{48} \nu^3 r^3 + O(r^4)$$

Noting the differing signs, we find:

$$r^* = \left(\frac{4}{20 + 3n_L} \right) \frac{1}{\nu}.$$

Optimal Aspect Ratio

Defining for the four-point vertex ratio,

$$\frac{V^{(L)}}{n_{L-1} (G^{(L)})^2} \equiv \nu r,$$

with ν an *activation function-dependent constant*, we get:

$$\mathcal{I}[p(\mathcal{M}_1, \mathcal{M}_2|x)] = \frac{m_1 m_2}{8} \nu^2 r^2 - \frac{m_1 m_2 (20 + 3m_1 + 3m_2)}{48} \nu^3 r^3 + O(r^4)$$

Noting the differing signs, we find:

$$r^* = \left(\frac{4}{20 + 3n_L} \right) \frac{1}{\nu}.$$

- For $n_L = 10$, $r^* = .12$ (tanh) and $r^* = .016$ (ReLU).

Optimal Aspect Ratio

Defining for the four-point vertex ratio,

$$\frac{V^{(L)}}{n_{L-1} (G^{(L)})^2} \equiv \nu r,$$

with ν an *activation function-dependent constant*, we get:

$$\mathcal{I}[p(\mathcal{M}_1, \mathcal{M}_2|x)] = \frac{m_1 m_2}{8} \nu^2 r^2 - \frac{m_1 m_2 (20 + 3m_1 + 3m_2)}{48} \nu^3 r^3 + O(r^4)$$

Noting the differing signs, we find:

$$r^* = \left(\frac{4}{20 + 3n_L} \right) \frac{1}{\nu}.$$

- ▶ For $n_L = 10$, $r^* = .12$ (tanh) and $r^* = .016$ (ReLU).
- ▶ **Residual connections** let us push r^* to arbitrary depths.

Epilogue: Model Complexity

According to the hype of 1987, neural networks were meant to be intelligent models that discovered features and patterns in data. Gaussian processes in contrast are simply smoothing devices. How can Gaussian processes possibly replace neural networks? Were neural networks over-hyped, or have we underestimated the power of smoothing methods?

David MacKay

- ▶ The success of **overparameterized** models with far more parameters than training data has led many to conjecture that “more is better” when it comes to deep learning.
- ▶ There's mounting *empirical* evidence that a **scaling hypothesis** captures the behavior of deep neural networks, signaling the optimality of the overparameterized regime.

Epilogue: Model Complexity

The **Occam's razor principle of sparsity** posits that we should favor the simplest hypothesis that explains our observations:

- ▶ In the context of machine learning, this is usually interpreted to mean that we should prefer models with fewer parameters when comparing models performing the same tasks.
- ▶ We expect that models with fewer parameters will have smaller *generalization errors*, and will *overfit* less.

Epilogue: Model Complexity

The **Occam's razor principle of sparsity** posits that we should favor the simplest hypothesis that explains our observations:

- ▶ In the context of machine learning, this is usually interpreted to mean that we should prefer models with fewer parameters when comparing models performing the same tasks.
- ▶ We expect that models with fewer parameters will have smaller *generalization errors*, and will *overfit* less.

To conclude our lectures, we're going to see how to resolve this puzzling within our framework.

Epilogue: Model Complexity

This hinges on the notion of **model complexity**:

Epilogue: Model Complexity

This hinges on the notion of **model complexity**:

- ▶ On the one hand, the orthodox discussion of generalization takes a **microscopic perspective** – focusing on how a network works in terms of its low-level components – and wants to identify model complexity with *model parameters*.

Epilogue: Model Complexity

This hinges on the notion of **model complexity**:

- ▶ On the one hand, the orthodox discussion of generalization takes a **microscopic perspective** – focusing on how a network works in terms of its low-level components – and wants to identify model complexity with *model parameters*.
- ▶ On the other hand, in these lectures we integrated out the model parameters and developed a **macroscopic perspective** – providing an effective theory description of the predictions of realistic fully-trained networks – for which this notion of model complexity is completely *reversed*.

Epilogue: Model Complexity

We now know it's the **depth-to-width aspect ratio**,

$$r \equiv L/n,$$

controlling the complexity of overparameterized neural networks.

- ▶ It's the number of **data-dependent couplings** specifying the truncated nearly-Gaussian distribution – and *not* the number of *model parameters* – that ultimately define the model complexity in deep learning.

Sparsity at Infinite Width

At *infinite width*, we found a **Gaussian** trained distribution:

$$\lim_{n \rightarrow \infty} p(z(\infty)) \equiv p(z(\infty) | y_{\tilde{\alpha}}, K_{\delta_1 \delta_2}, \Theta_{\delta_1 \delta_2}) .$$

Sparsity at Infinite Width

At *infinite width*, we found a **Gaussian** trained distribution:

$$\lim_{n \rightarrow \infty} p(z(\infty)) \equiv p(z(\infty) | y_{\tilde{\alpha}}, K_{\delta_1 \delta_2}, \Theta_{\delta_1 \delta_2}) .$$

- ▶ The reason for writing it as a *conditional distribution* in this way is that the mean is only a function of $y_{\tilde{\alpha}}$ and $\Theta_{\delta_1 \delta_2}^{(L)}$, while the variance is only a function of $K_{\delta_1 \delta_2}^{(L)}$ and $\Theta_{\delta_1 \delta_2}^{(L)}$.

Sparsity at Infinite Width

At *infinite width*, we found a **Gaussian** trained distribution:

$$\lim_{n \rightarrow \infty} p(z(\infty)) \equiv p(z(\infty) | y_{\tilde{\alpha}}, K_{\delta_1 \delta_2}, \Theta_{\delta_1 \delta_2}) .$$

- ▶ The reason for writing it as a *conditional distribution* in this way is that the mean is only a function of $y_{\tilde{\alpha}}$ and $\Theta_{\delta_1 \delta_2}^{(L)}$, while the variance is only a function of $K_{\delta_1 \delta_2}^{(L)}$ and $\Theta_{\delta_1 \delta_2}^{(L)}$.
- ▶ This is **sparse**, depending on a few objects in a simple way.

Near-Sparsity at Finite Width

At *finite width*, we found a **nearly-Gaussian** trained distribution:

$$p(z(\infty)) \equiv p(z(\infty) | y, G, H, V, A, B, D, F, P, Q, R, S, T, U) + O\left(\frac{1}{n^2}\right).$$

Near-Sparsity at Finite Width

At *finite width*, we found a **nearly-Gaussian** trained distribution:

$$p(z(\infty)) \equiv p(z(\infty) | y, G, H, V, A, B, D, F, P, Q, R, S, T, U) + O\left(\frac{1}{n^2}\right).$$

- ▶ In addition to G and H , here we're accounting for the finite-width *data-dependent couplings* arising from:

$$\begin{aligned} & \mathbb{E} [zzzz]_{\text{connected}}, & \mathbb{E} [\widehat{\Delta H}zz], & \mathbb{E} [\widehat{\Delta H}^2], \\ & \mathbb{E} [\widehat{dHz}], & \mathbb{E} [\widehat{dd}_I H], & \mathbb{E} [\widehat{dd}_{II} H]. \end{aligned}$$

Near-Sparsity at Finite Width

At *finite width*, we found a **nearly-Gaussian** trained distribution:

$$p(z(\infty)) \equiv p(z(\infty) | y, G, H, V, A, B, D, F, P, Q, R, S, T, U) + O\left(\frac{1}{n^2}\right).$$

- ▶ In addition to G and H , here we're accounting for the finite-width *data-dependent couplings* arising from:

$$\begin{aligned} & \mathbb{E} [zzzz]_{\text{connected}}, & \mathbb{E} [\widehat{\Delta H}zz], & \mathbb{E} [\widehat{\Delta H}^2], \\ & \mathbb{E} [\widehat{dHz}], & \mathbb{E} [\widehat{dd}_I H], & \mathbb{E} [\widehat{dd}_{II} H]. \end{aligned}$$

- ▶ This is **nearly-sparse**, depending only on two-hands-full of objects in a *nearly-simple* way.

Model Complexity of Fully-Trained Neural Networks

Consider a fixed combined training and test dataset of size $N_{\mathcal{D}}$:

- ▶ For the *infinite-width* **Gaussian distribution**, we only need

$$n_{\text{out}}N_{\mathcal{A}} + \left[\frac{N_{\mathcal{D}}(N_{\mathcal{D}} + 1)}{2} \right] + \left[\frac{N_{\mathcal{D}}(N_{\mathcal{D}} + 1)}{2} \right] = O(N_{\mathcal{D}}^2)$$

numbers in order to completely specify the distribution.

Model Complexity of Fully-Trained Neural Networks

Consider a fixed combined training and test dataset of size $N_{\mathcal{D}}$:

- ▶ For the *finite-width* **nearly-Gaussian distribution** with $0 < r \ll 1$, we will instead need $O(N_{\mathcal{D}}^4)$ numbers, with the counting dominated by the finite-width tensors.

Model Complexity of Fully-Trained Neural Networks

Consider a fixed combined training and test dataset of size N_D :

- ▶ For an accuracy $O(L^k/n^k)$, a *macroscopic description*

$$p(z(\infty)) = \sum_{m=0}^k \frac{p^{\{m\}}(z(\infty))}{n^m} + O\left(\frac{L^{k+1}}{n^{k+1}}\right),$$

will need $O(N_D^{2k})$ numbers in general.

Model Complexity of Fully-Trained Neural Networks

Consider a fixed combined training and test dataset of size N_D :

- ▶ For an accuracy $O(L^k/n^k)$, a *macroscopic description*

$$p(z(\infty)) = \sum_{m=0}^k \frac{p^{\{m\}}(z(\infty))}{n^m} + O\left(\frac{L^{k+1}}{n^{k+1}}\right),$$

will need $O(N_D^{2k})$ numbers in general.

- ▶ The **1/n expansion** gives a sequence of effective theories with increasing accuracy at the cost of increasing complexity.

Model Complexity of Fully-Trained Neural Networks

What we have found here is the manifestation of the **microscopic-macroscopic duality**:

Model Complexity of Fully-Trained Neural Networks

What we have found here is the manifestation of the **microscopic-macroscopic duality**:

- ▶ Complexity in *parameter space* is traded into simplicity in *sample space*, and density in *model parameters* is exchanged for sparsity in *data-dependent couplings*.

Model Complexity of Fully-Trained Neural Networks

What we have found here is the manifestation of the **microscopic-macroscopic duality**:

- ▶ Complexity in *parameter space* is traded into simplicity in *sample space*, and density in *model parameters* is exchanged for sparsity in *data-dependent couplings*.
- ▶ In the *overparameterized regime*, this indicates that we should identify the *model complexity* with the **data-dependent couplings** rather than the **model parameters**.

Model Complexity of Fully-Trained Neural Networks

As r increases, we'll need to include more of these higher-order terms, making our macroscopic description more complex:

Model Complexity of Fully-Trained Neural Networks

As r increases, we'll need to include more of these higher-order terms, making our macroscopic description more complex:

- ▶ In the strict limit $r \rightarrow 0$, the *sparse* $O(N_D^2)$ **Gaussian** description of the infinite-width limit will be accurate.

Model Complexity of Fully-Trained Neural Networks

As r increases, we'll need to include more of these higher-order terms, making our macroscopic description more complex:

- ▶ In the strict limit $r \rightarrow 0$, the *sparse* $O(N_{\mathcal{D}}^2)$ **Gaussian** description of the infinite-width limit will be accurate.
- ▶ In the regime $0 < r \sim r^* \ll 1$, the *nearly-sparse* $O(N_{\mathcal{D}}^4)$ **nearly-Gaussian** description of the finite-width effective theory truncated at order $1/n$ will be accurate.

Model Complexity of Fully-Trained Neural Networks

As r increases, we'll need to include more of these higher-order terms, making our macroscopic description more complex:

- ▶ In the strict limit $r \rightarrow 0$, the *sparse* $O(N_{\mathcal{D}}^2)$ **Gaussian** description of the infinite-width limit will be accurate.
- ▶ In the regime $0 < r \sim r^* \ll 1$, the *nearly-sparse* $O(N_{\mathcal{D}}^4)$ **nearly-Gaussian** description of the finite-width effective theory truncated at order $1/n$ will be accurate.
- ▶ For larger r , a more generic $O(N_{\mathcal{D}}^{2k})$ **non-Gaussian** description would in principle be necessary.

Conclusion

The practical success of deep learning in the *overparameterized* regime and the empirical accuracy of a simple *scaling hypothesis* is really telling us that useful neural networks should be **sparse** – hence the preference for larger and larger models – but not too sparse – so that they are also **deep**. Thus, from the macroscopic perspective, a **nearly-sparse** model complexity is perhaps the most important inductive bias of deep learning.

Conclusion

The practical success of deep learning in the *overparameterized* regime and the empirical accuracy of a simple *scaling hypothesis* is really telling us that useful neural networks should be **sparse** – hence the preference for larger and larger models – but not too sparse – so that they are also **deep**. Thus, from the macroscopic perspective, a **nearly-sparse** model complexity is perhaps the most important inductive bias of deep learning.

Thank You!