

Non-Gaussian processes and neural networks at finite widths

Sho Yaida

Facebook AI Research

Facebook Inc.

Menlo Park, California 94025, USA

SHOYaida@fb.com

Abstract

Gaussian processes are ubiquitous in nature and engineering. A case in point is a class of neural networks in the infinite-width limit, whose priors correspond to Gaussian processes. Here we perturbatively extend this correspondence to finite-width neural networks, yielding non-Gaussian processes as priors. The methodology developed herein allows us to track the flow of preactivation distributions by progressively integrating out random variables from lower to higher layers, reminiscent of renormalization-group flow. We further develop a perturbative procedure to perform Bayesian inference with weakly non-Gaussian priors.

1. Inception

Gaussian processes model many phenomena in the physical world. A prime example is Brownian motion (Brown, 1828), modeled as the integral of Gaussian-distributed bumps exerted on a point-like solute (Einstein, 1905). The theory of elementary particles (Weinberg, 1995) also becomes a Gaussian process in the free limit where interactions between particles are turned off, and many-body systems as complex as glasses come to be Gaussian in the infinite-dimensional, mean-field, limit (Parisi and Zamponi, 2010). In the context of machine learning, Neal (1996) pointed out that a class of neural networks give rise to Gaussian processes in the infinite-width limit, which can perform exact Bayesian inference from training to test data (Williams, 1997). They occupy a corner of theoretical playground wherein the *karakuri* of neural networks is scrutinized (Lee et al., 2018; Matthews et al., 2018; Jacot et al., 2018; Chizat et al., 2018; Lee et al., 2019; Geiger et al., 2019).

In reality, Gaussian processes are but mere idealizations. Brownian particles have finite-size structure, elementary particles interact, and many-body systems respond nonlinearly. In order to understand rich phenomena exhibited by these real systems, Gaussian processes rather serve as starting points to be perturbed around. Indeed many edifices in theoretical physics are built upon the successful treatment of non-Gaussianity, with a notable example being renormalization-group flow (Kadanoff, 1966; Wilson, 1971; Weinberg, 1996; Goldenfeld, 2018). In the quest to elucidate behaviors of real neural networks away from the infinite-width limit, it is thus natural to wonder if the similar treatment of non-Gaussianity yields equally elegant and powerful machinery.

Here we set out on this program, perturbatively treating finite-width corrections to neural networks. Prior distributions of outputs are obtained through progressively integrating out preactivation of neurons layer by layer, yielding non-Gaussian priors. The whole procedure closely resembles renormalization-group flow (Goldenfeld, 2018; Mehta and Schwab, 2014): it bridges probability distributions at different scales through coarse-graining of random variables at microscopic scales; the flow of distributions is traced through running couplings, which in particular capture the de-

gree of non-Gaussianity in these distributions; resulting recursive equations (R1,R2,R3) govern the evolution of these running couplings from lower to higher layers, just as renormalization-group equations do from microscopic to macroscopic scales. Such a recursive approach enables us to treat finite-width corrections to various observables, for networks with arbitrary activation functions.

The rest of the paper is structured as follows. In Section 2 we review and set up basic concepts. Our master recursive formulae (R1,R2,R3) are derived in Section 3, which control the flow of preactivation distributions. After an interlude with concrete examples in Section 4, we extend the Gaussian-process Bayesian inference to non-Gaussian priors in Section 5 and study inference of neural networks at finite widths. We conclude in Section 6 with dreams.

2. To infinity and beyond

In this paper we study real finite-width neural networks in the regime where the number of neurons in hidden layers is asymptotically large whereas input and output dimensions are kept constant.

2.1. Gaussian processes and neural networks at infinite widths

Let us focus on a class of neural networks termed multilayer perceptrons, with model parameters, $\theta = \{b_i^{(\ell)}, W_{ij}^{(\ell)}\}$, and an activation function, σ . For each input, $\mathbf{x} \in \mathbb{R}^{n_0}$, a neural network outputs a vector, $\mathbf{z}(\mathbf{x}; \theta) = \mathbf{z}^{(L)} \in \mathbb{R}^{n_L}$, recursively defined as sequences of preactivations through

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

$$z_i^{(\ell)}(\mathbf{x}) = b_i^{(\ell)} + \sum_{j=1}^{n_{\ell-1}} W_{ij}^{(\ell)} \sigma \left[z_j^{(\ell-1)}(\mathbf{x}) \right] \quad \text{for } i = 1, \dots, n_\ell; \quad \ell = 2, \dots, L. \quad (2)$$

Following Neal (1996), we assume priors for biases and weights given by independent and identically distributed Gaussian distributions with zero means, $\mathbb{E} \left[b_i^{(\ell)} \right] = \mathbb{E} \left[W_{ij}^{(\ell)} \right] = 0$, and variances

$$\mathbb{E} \left[b_{i_1}^{(\ell)} b_{i_2}^{(\ell)} \right] = \delta_{i_1 i_2} C_b^{(\ell)}, \quad (3)$$

$$\mathbb{E} \left[W_{i_1 j_1}^{(\ell)} W_{i_2 j_2}^{(\ell)} \right] = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{C_W^{(\ell)}}{n_{\ell-1}}. \quad (4)$$

Higher moments are then obtained by Wick's contractions (Wick, 1950; Zee, 2010). For instance,

$$\mathbb{E} \left[b_{i_1}^{(\ell)} b_{i_2}^{(\ell)} b_{i_3}^{(\ell)} b_{i_4}^{(\ell)} \right] = \left[C_b^{(\ell)} \right]^2 \times (\delta_{i_1 i_2} \delta_{i_3 i_4} + \delta_{i_1 i_3} \delta_{i_2 i_4} + \delta_{i_1 i_4} \delta_{i_2 i_3}). \quad (5)$$

For those unfamiliar with Wick's contractions and connected correlation functions (a.k.a. cumulants), a pedagogical review is provided in Appendix A as our formalism heavily relies on them.

In the infinite-width limit where $n_1, n_2, \dots, n_{L-1} \rightarrow \infty$ (but finite n_0 and n_L), it has been argued – with varying degrees of rigor (Neal, 1996; Lee et al., 2018; Matthews et al., 2018) – that the prior distribution of outputs is governed by the Gaussian process with a kernel

$$K_{i_1 i_2; \alpha_1 \alpha_2} \equiv \mathbb{E} \left[z_{i_1}^{(L)}(\mathbf{x}_{\alpha_1}) z_{i_2}^{(L)}(\mathbf{x}_{\alpha_2}) \right] \quad (6)$$

and all the higher moments given by Wick’s contractions. Here, the sample index α labels different inputs in a dataset. There exists a recursive formula that lets us evaluate this kernel for any pair of inputs (Lee et al., 2018) [c.f. Equation (R1)]. Importantly, once the values of the kernel are evaluated for all the pairs of $N_D = N_R + N_E$ input data, $\{\mathbf{x}_\alpha\}_{\alpha=1,\dots,N_D}$, consisting of N_R training inputs with target outputs and N_E test inputs with unknown targets, we can perform exact Bayesian inference to yield mean outputs as predictions for N_E test data (Williams, 1997; Williams and Rasmussen, 2006) [c.f. Equation (GPM)]. This should be contrasted with stochastic gradient descent (SGD) optimization (Robbins and Monro, 1951), through which typically a single estimate for the optimal model parameters of the posterior, θ_* , is obtained and used to predict outputs for test inputs; Bayesian inference instead marginalizes over all model parameters, performing an ensemble average over the posterior distribution (MacKay, 1995).

2.2. Beyond infinity

We shall now study real finite-width neural networks in the regime $n_1, \dots, n_{L-1} \sim n \gg 1$.¹ At finite widths, there are corrections to Gaussian-process priors. In other words, a whole tower of nontrivial preactivation correlation functions beyond the kernel,

$$G_{i_1 \dots i_m; \alpha_1 \dots \alpha_m}^{(\ell)} \equiv \mathbb{E} \left[z_{i_1}^{(\ell)}(\mathbf{x}_{\alpha_1}) \cdots z_{i_m}^{(\ell)}(\mathbf{x}_{\alpha_m}) \right], \quad (7)$$

collectively dictate the distribution of preactivations. Our aim is to trace the flow of these distributions progressively and cumulatively all the way up to the last layer whereat Bayesian inference is executed. More specifically, we shall inductively and self-consistently show that two-point preactivation correlation functions take the form²

$$G_{i_1 i_2; \alpha_1 \alpha_2}^{(\ell)} = \delta_{i_1 i_2} \left[\tilde{K}_{\alpha_1 \alpha_2}^{(\ell)} + \frac{1}{n_{\ell-1}} \tilde{S}_{\alpha_1 \alpha_2}^{(\ell)} + O\left(\frac{1}{n^2}\right) \right] \quad (\text{KS})$$

and connected four-point preactivation correlation functions

$$\begin{aligned} & G_{i_1 i_2 i_3 i_4; \alpha_1 \alpha_2 \alpha_3 \alpha_4}^{(\ell)} \Big|_{\text{connected}} \quad (V) \\ & \equiv G_{i_1 i_2 i_3 i_4; \alpha_1 \alpha_2 \alpha_3 \alpha_4}^{(\ell)} - G_{i_1 i_2; \alpha_1 \alpha_2}^{(\ell)} G_{i_3 i_4; \alpha_3 \alpha_4}^{(\ell)} - G_{i_1 i_3; \alpha_1 \alpha_3}^{(\ell)} G_{i_2 i_4; \alpha_2 \alpha_4}^{(\ell)} - G_{i_1 i_4; \alpha_1 \alpha_4}^{(\ell)} G_{i_2 i_3; \alpha_2 \alpha_3}^{(\ell)} \\ & = \frac{1}{n_{\ell-1}} \left[\delta_{i_1 i_2} \delta_{i_3 i_4} \tilde{V}_{(\alpha_1 \alpha_2)(\alpha_3 \alpha_4)}^{(\ell)} + \delta_{i_1 i_3} \delta_{i_2 i_4} \tilde{V}_{(\alpha_1 \alpha_3)(\alpha_2 \alpha_4)}^{(\ell)} + \delta_{i_1 i_4} \delta_{i_2 i_3} \tilde{V}_{(\alpha_1 \alpha_4)(\alpha_2 \alpha_3)}^{(\ell)} \right] + O\left(\frac{1}{n^2}\right), \end{aligned}$$

and higher cumulants are all suppressed by $O\left(\frac{1}{n^2}\right)$.³ Here the Gaussian-process core kernel $\tilde{K}_{\alpha_1 \alpha_2}^{(\ell)}$ and the self-energy correction $\tilde{S}_{\alpha_1 \alpha_2}^{(\ell)}$ are symmetric under the exchange of sample indices $\alpha_1 \leftrightarrow \alpha_2$ and the four-point vertex $\tilde{V}_{(\alpha_1 \alpha_2)(\alpha_3 \alpha_4)}^{(\ell)}$ is symmetric under $\alpha_1 \leftrightarrow \alpha_2$, $\alpha_3 \leftrightarrow \alpha_4$, and $(\alpha_1 \alpha_2) \leftrightarrow (\alpha_3 \alpha_4)$.

1. Note that input and output dimensions, n_0 and n_L , are arbitrary. To be precise, defining $n_1, \dots, n_{L-1} \equiv r_1 n, \dots, r_{L-1} n$, we send $n \gg 1$ while keeping $\left\{ C_b^{(\ell)}, C_W^{(\ell)} \right\}_{\ell=1, \dots, L}$, r_1, \dots, r_{L-1} , n_0 , and n_L constants, and compute the leading $1/n$ corrections. In particular it is crucial to keep the number of outputs n_L constant in order to consistently perform Bayesian inference within our approach.
2. In the main text we place tildes on objects that depend only on sample indices α ’s in order to distinguish them from those that depend both on sample indices α ’s and neuron indices i ’s.
3. Given that the means of biases and weights are zero, $G_{i_1 \dots i_m; \alpha_1 \dots \alpha_m}^{(\ell)} = 0$ for all odd m .

$(\alpha_3\alpha_4)$. At the first layer the preactivation distribution is exactly Gaussian for any finite widths and hence Equations (KS) and (V) are trivially satisfied, with

$$\tilde{K}_{\alpha_1\alpha_2}^{(1)} = C_b^{(1)} + C_W^{(1)} \cdot \left(\frac{\mathbf{x}_{\alpha_1} \cdot \mathbf{x}_{\alpha_2}}{n_0} \right), \quad \tilde{S}_{\alpha_1\alpha_2}^{(1)} = 0, \quad \text{and} \quad \tilde{V}_{(\alpha_1\alpha_2)(\alpha_3\alpha_4)}^{(1)} = 0. \quad (\text{R0})$$

Obtained in Section 3 are the recursive formulae that link these core kernel, self-energy, and four-point vertex at the ℓ -th layer to those at the $(\ell + 1)$ -th layer while in Section 5 these tensors at the last layer $\ell = L$ are used to yield the leading $1/n$ correction for Bayesian inference at finite widths.

2.3. Related work

Our Schwinger operator approach is orthogonal to the replica approach by Cohen et al. (2019) and, unlike the planar diagrammatic approach by Dyer and Gur-Ari (2019), applies to general activation functions, made possible by accumulating corrections layer by layer rather than dealing with them all at once. See also Antognini (2019). More substantially, in contrast to these previous approaches, we here study finite-width effects on Bayesian inference and find that the renormalization-group picture naturally emerges, with layers playing the role of scales.

3. Distributional flow

As auxiliary objects in recursive steps, let us introduce activation correlation functions

$$H_{i_1 \dots i_m; \alpha_1 \dots \alpha_m}^{(\ell)} \equiv \mathbb{E} \left\{ \sigma \left[z_{i_1}^{(\ell)}(\mathbf{x}_{\alpha_1}) \right] \cdots \sigma \left[z_{i_m}^{(\ell)}(\mathbf{x}_{\alpha_m}) \right] \right\}. \quad (8)$$

Our basic strategy is to establish relations

$$\left\{ \mathbf{G}^{(1)} \right\} \rightarrow \left\{ \mathbf{H}^{(1)} \right\} \rightarrow \left\{ \mathbf{G}^{(2)} \right\} \rightarrow \cdots \rightarrow \left\{ \mathbf{H}^{(L-1)} \right\} \rightarrow \left\{ \mathbf{G}^{(L)} \right\}, \quad (\text{ZIGZAG})$$

zigzagging between sets of preactivation correlation functions and sets of activation correlation functions, keeping track of leading finite-width corrections. Below, relations $\mathbf{G}^{(\ell)} \rightarrow \mathbf{H}^{(\ell)}$ are obtained by integrating out preactivations while relations $\mathbf{H}^{(\ell)} \rightarrow \mathbf{G}^{(\ell+1)}$ are obtained by integrating out biases and weights. At first glance the algebra in this paper may look horrifying but repeated applications of Wick's contractions are all there is to it. The results are summarized in Section 3.2.

3.1. Zigzag relations for preactivation and activation correlation functions

Integrating over the Gaussian biases and weights at ℓ 's connections yield the relations that link activation correlations $\mathbf{H}^{(\ell)}$ to preactivation correlations $\mathbf{G}^{(\ell+1)}$ at the next layer. Recalling Equations (KS) and (V), trivial Wick's contractions yield

$$\tilde{K}_{\alpha_1\alpha_2}^{(\ell+1)} + \frac{1}{n_\ell} \tilde{S}_{\alpha_1\alpha_2}^{(\ell+1)} = C_b^{(\ell+1)} + C_W^{(\ell+1)} \left[\frac{1}{n_\ell} \sum_{j=1}^{n_\ell} H_{jj; \alpha_1\alpha_2}^{(\ell)} \right] + O\left(\frac{1}{n^2}\right) \quad \text{and} \quad (9)$$

$$\tilde{V}_{(\alpha_1\alpha_2)(\alpha_3\alpha_4)}^{(\ell+1)} = \frac{[C_W^{(\ell+1)}]^2}{n_\ell} \sum_{j,k=1}^{n_\ell} \left[H_{jjkk; \alpha_1\alpha_2\alpha_3\alpha_4}^{(\ell)} - H_{jj; \alpha_1\alpha_2}^{(\ell)} H_{kk; \alpha_3\alpha_4}^{(\ell)} \right] + O\left(\frac{1}{n}\right). \quad (10)$$

The remaining task is to relate preactivation correlations $\mathbf{G}^{(\ell)}$ to activation correlations $\mathbf{H}^{(\ell)}$ within the same layer, which will complete the zigzag relation (**ZIGZAG**) for these correlation functions.⁴

With the mastery of Wick's contractions and connected correlation functions, it is simple to derive the following combinatorial hack (Appendix A.4): viewing prior preactivations

$$\mathbf{z} \equiv \left\{ \mathbf{z}_{i;\alpha} \equiv z_i^{(\ell)}(\mathbf{x}_\alpha) \right\}_{i=1,\dots,n_\ell; \alpha=1,\dots,N_D}$$

at the ℓ -th layer as a random $(n_\ell N_D)$ -dimensional vector and defining the Gaussian integral with the kernel $\langle \mathbf{z}_{i_1;\alpha_1} \mathbf{z}_{i_2;\alpha_2} \rangle_{K^{(\ell)}} = K_{i_1 i_2; \alpha_1 \alpha_2}^{(\ell)} \equiv \delta_{i_1 i_2} \tilde{K}_{\alpha_1 \alpha_2}^{(\ell)}$, the prior average

$$\mathbb{E} \{ \mathcal{F}[\mathbf{z}] \} = \langle \mathcal{F}[\mathbf{z}] \rangle_{K^{(\ell)}} + \frac{1}{n_{\ell-1}} [\langle \mathcal{F}[\mathbf{z}] \mathcal{O}_S[\mathbf{z}] + \mathcal{F}[\mathbf{z}] \mathcal{O}_V[\mathbf{z}] \rangle_{K^{(\ell)}}] + O\left(\frac{1}{n^2}\right) \quad (\text{HACK})$$

for any function \mathcal{F} . Here the operators $\mathcal{O}_S[\mathbf{z}]$ and $\mathcal{O}_V[\mathbf{z}]$ capture $1/n$ corrections due to self-energy and four-point vertex, respectively, and are defined as

$$\mathcal{O}_S[\mathbf{z}] \equiv \frac{1}{2} \sum_{\alpha_1, \alpha_2} \tilde{S}_{(\ell)}^{\alpha_1 \alpha_2} \left[\left(\sum_{i=1}^{n_\ell} \mathbf{z}_{i;\alpha_1} \mathbf{z}_{i;\alpha_2} \right) - n_\ell \tilde{K}_{\alpha_1 \alpha_2}^{(\ell)} \right] \quad \text{and} \quad (\text{OS})$$

$$\begin{aligned} \mathcal{O}_V[\mathbf{z}] &\equiv \frac{1}{8} \sum_{\alpha_1, \alpha_2, \alpha_3, \alpha_4} \tilde{V}_{(\ell)}^{(\alpha_1 \alpha_2)(\alpha_3 \alpha_4)} \\ &\times \left[\left(\sum_{i=1}^{n_\ell} \mathbf{z}_{i;\alpha_1} \mathbf{z}_{i;\alpha_2} \right) \left(\sum_{j=1}^{n_\ell} \mathbf{z}_{j;\alpha_3} \mathbf{z}_{j;\alpha_4} \right) - 2n_\ell \left(\sum_{i=1}^{n_\ell} \mathbf{z}_{i;\alpha_1} \mathbf{z}_{i;\alpha_2} \right) \tilde{K}_{\alpha_3 \alpha_4}^{(\ell)} \right. \\ &\left. - 4 \left(\sum_{i=1}^{n_\ell} \mathbf{z}_{i;\alpha_1} \mathbf{z}_{i;\alpha_3} \right) \tilde{K}_{\alpha_2 \alpha_4}^{(\ell)} + n_\ell^2 \tilde{K}_{\alpha_1 \alpha_2}^{(\ell)} \tilde{K}_{\alpha_3 \alpha_4}^{(\ell)} + 2n_\ell \tilde{K}_{\alpha_1 \alpha_3}^{(\ell)} \tilde{K}_{\alpha_2 \alpha_4}^{(\ell)} \right], \quad (\text{OV}) \end{aligned}$$

where the sample indices are raised by using the inverse core kernel as a metric, meaning

$$\tilde{S}_{(\ell)}^{\alpha_1 \alpha_2} \equiv \sum_{\alpha'_1, \alpha'_2} \left(\tilde{K}_{(\ell)}^{-1} \right)^{\alpha_1 \alpha'_1} \left(\tilde{K}_{(\ell)}^{-1} \right)^{\alpha_2 \alpha'_2} \tilde{S}_{\alpha'_1 \alpha'_2}^{(\ell)} \quad \text{and} \quad (11)$$

$$\tilde{V}_{(\ell)}^{(\alpha_1 \alpha_2)(\alpha_3 \alpha_4)} \equiv \sum_{\alpha'_1, \dots, \alpha'_4} \left(\tilde{K}_{(\ell)}^{-1} \right)^{\alpha_1 \alpha'_1} \dots \left(\tilde{K}_{(\ell)}^{-1} \right)^{\alpha_4 \alpha'_4} \tilde{V}_{(\alpha'_1 \alpha'_2)(\alpha'_3 \alpha'_4)}^{(\ell)}. \quad (12)$$

Using the above hack, we can evaluate the activation correlations by straightforward algebra with Wick's contractions. In particular, as the Gaussian integral is diagonal in the neuron index i , we just need to disentangle cases with repeated and unrepeated neuron indices. The solution for this exercise is in Appendix B: it is arguably the most cumbersome algebra in this paper.

4. The nontrivial parts of the inductive proof for Equations (KS) and (V) are to show (i) that the right-hand side of Equation (10) is finite as $n \rightarrow \infty$, (ii) that the leading contribution of Equation (9) is the Gaussian-process kernel, and (iii) that higher-point connected preactivation correlation functions are all suppressed by $O\left(\frac{1}{n^2}\right)$, all of which are verified in obtaining the recursive equations. See Appendix B for a full proof.

3.2. Master recursive flow equations

Denoting the Gaussian integral with the core kernel $\langle \tilde{\mathbf{z}}_{\alpha_1} \tilde{\mathbf{z}}_{\alpha_2} \rangle_{\tilde{K}^{(\ell)}} = \tilde{K}_{\alpha_1 \alpha_2}^{(\ell)}$ for a single-neuron random vector $\tilde{\mathbf{z}} \equiv \{\tilde{z}_\alpha\}_{\alpha=1, \dots, N_D}$, and plugging in results of Appendix B into Equations (9) and (10), we arrive at our master recursion relations

$$\tilde{K}_{\alpha_1 \alpha_2}^{(\ell+1)} = C_b^{(\ell+1)} + C_W^{(\ell+1)} \langle \sigma(\tilde{\mathbf{z}}_{\alpha_1}) \sigma(\tilde{\mathbf{z}}_{\alpha_2}) \rangle_{\tilde{K}^{(\ell)}}, \quad (\text{R1})$$

$$\begin{aligned} \tilde{V}_{(\alpha_1 \alpha_2)(\alpha_3 \alpha_4)}^{(\ell+1)} = & \left[C_W^{(\ell+1)} \right]^2 \left[\langle \sigma(\tilde{\mathbf{z}}_{\alpha_1}) \sigma(\tilde{\mathbf{z}}_{\alpha_2}) \sigma(\tilde{\mathbf{z}}_{\alpha_3}) \sigma(\tilde{\mathbf{z}}_{\alpha_4}) \rangle_{\tilde{K}^{(\ell)}} \right. \\ & - \langle \sigma(\tilde{\mathbf{z}}_{\alpha_1}) \sigma(\tilde{\mathbf{z}}_{\alpha_2}) \rangle_{\tilde{K}^{(\ell)}} \langle \sigma(\tilde{\mathbf{z}}_{\alpha_3}) \sigma(\tilde{\mathbf{z}}_{\alpha_4}) \rangle_{\tilde{K}^{(\ell)}} \\ & + \frac{1}{4} \left(\frac{n_\ell}{n_{\ell-1}} \right) \sum_{\alpha'_1, \alpha'_2, \alpha'_3, \alpha'_4} \tilde{V}_{(\alpha'_1 \alpha'_2)(\alpha'_3 \alpha'_4)}^{(\ell)} \left\langle \sigma(\tilde{\mathbf{z}}_{\alpha_1}) \sigma(\tilde{\mathbf{z}}_{\alpha_2}) (\tilde{z}_{\alpha'_1} \tilde{z}_{\alpha'_2} - \tilde{K}_{\alpha'_1 \alpha'_2}^{(\ell)}) \right\rangle_{\tilde{K}^{(\ell)}} \\ & \left. \times \left\langle \sigma(\tilde{\mathbf{z}}_{\alpha_3}) \sigma(\tilde{\mathbf{z}}_{\alpha_4}) (\tilde{z}_{\alpha'_3} \tilde{z}_{\alpha'_4} - \tilde{K}_{\alpha'_3 \alpha'_4}^{(\ell)}) \right\rangle_{\tilde{K}^{(\ell)}} \right], \quad \text{and} \end{aligned} \quad (\text{R2})$$

$$\begin{aligned} \tilde{S}_{\alpha_1 \alpha_2}^{(\ell+1)} = & \left(\frac{n_\ell}{n_{\ell-1}} \right) C_W^{(\ell+1)} \left[\frac{1}{2} \sum_{\alpha'_1, \alpha'_2} \tilde{S}_{\alpha'_1 \alpha'_2}^{(\ell)} \left\langle \sigma(\tilde{\mathbf{z}}_{\alpha_1}) \sigma(\tilde{\mathbf{z}}_{\alpha_2}) (\tilde{z}_{\alpha'_1} \tilde{z}_{\alpha'_2} - \tilde{K}_{\alpha'_1 \alpha'_2}^{(\ell)}) \right\rangle_{\tilde{K}^{(\ell)}} \right. \\ & + \frac{1}{8} \sum_{\alpha'_1, \alpha'_2, \alpha'_3, \alpha'_4} \tilde{V}_{(\alpha'_1 \alpha'_2)(\alpha'_3 \alpha'_4)}^{(\ell)} \left\langle \sigma(\tilde{\mathbf{z}}_{\alpha_1}) \sigma(\tilde{\mathbf{z}}_{\alpha_2}) \right. \\ & \times \left(\tilde{z}_{\alpha'_1} \tilde{z}_{\alpha'_2} \tilde{z}_{\alpha'_3} \tilde{z}_{\alpha'_4} - 2\tilde{z}_{\alpha'_1} \tilde{z}_{\alpha'_2} \tilde{K}_{\alpha'_3 \alpha'_4}^{(\ell)} - 4\tilde{z}_{\alpha'_1} \tilde{z}_{\alpha'_3} \tilde{K}_{\alpha'_2 \alpha'_4}^{(\ell)} \right. \\ & \left. \left. + \tilde{K}_{\alpha'_1 \alpha'_2}^{(\ell)} \tilde{K}_{\alpha'_3 \alpha'_4}^{(\ell)} + 2\tilde{K}_{\alpha'_1 \alpha'_3}^{(\ell)} \tilde{K}_{\alpha'_2 \alpha'_4}^{(\ell)} \right) \right\rangle_{\tilde{K}^{(\ell)}} \left. \right]. \end{aligned} \quad (\text{R3})$$

For $\ell = 1$, a special note about the ratio $\frac{n_\ell}{n_{\ell-1}}$ is in order: even though n_0 stays constant while $n_1 \gg 1$, the terms proportional to that ratio are identically zero due to the complete Gaussianity (R0).

The preactivation distribution in the first layer (R0) sets the initial condition for the flow from lower to higher layers dictated by these recursive equations. Evolving through these recursive equations, the running couplings – $\tilde{K}_{\alpha_1 \alpha_2}^{(\ell)}$, $\tilde{V}_{(\alpha_1 \alpha_2)(\alpha_3 \alpha_4)}^{(\ell)}$, and $\tilde{S}_{\alpha_1 \alpha_2}^{(\ell)}$ – then trace changes in the distributions of preactivations as the layer scale ℓ shifts, just as running couplings for physical systems track changes in effective Boltzmann distributions as the probing scale shifts. Once recursed up to the last layer $\ell = L$, the resulting distribution of outputs $\mathbf{z} = \mathbf{z}^{(L)}$ can be succinctly encoded by the probability distribution

$$p[\mathbf{z}] = \frac{e^{-\mathcal{H}[\mathbf{z}]}}{\int d\mathbf{z}' e^{-\mathcal{H}[\mathbf{z}']}} \quad (\text{D0})$$

with the potential $\mathcal{H}[\mathbf{z}] = \mathcal{H}_0[\mathbf{z}] + \epsilon \mathcal{H}_1[\mathbf{z}] + O(\epsilon^2)$ where $\epsilon \equiv \frac{1}{n_{L-1}} \ll 1$,

$$\mathcal{H}_0[\mathbf{z}] = \frac{1}{2} \sum_{\alpha_1, \alpha_2} \left(\tilde{K}_{(L)}^{-1} \right)^{\alpha_1 \alpha_2} \left(\sum_{i=1}^{n_L} z_{i; \alpha_1} z_{i; \alpha_2} \right), \quad \text{and} \quad (\text{D1})$$

$$\mathcal{H}_1[\mathbf{z}] = -\frac{1}{2} \sum_{\alpha_1, \alpha_2} \tilde{J}^{\alpha_1 \alpha_2} \left(\sum_{i=1}^{n_L} z_{i; \alpha_1} z_{i; \alpha_2} \right) \quad (\text{D2})$$

$$-\frac{1}{8} \sum_{\alpha_1, \alpha_2, \alpha_3, \alpha_4} \tilde{V}_{(L)}^{(\alpha_1 \alpha_2)(\alpha_3 \alpha_4)} \left(\sum_{i=1}^{n_L} z_{i; \alpha_1} z_{i; \alpha_2} \right) \left(\sum_{j=1}^{n_L} z_{j; \alpha_3} z_{j; \alpha_4} \right) \quad \text{with} \\ \tilde{J}^{\alpha_1 \alpha_2} \equiv \tilde{S}_{(L)}^{\alpha_1 \alpha_2} - \sum_{\alpha_3, \alpha_4} \tilde{K}_{\alpha_3 \alpha_4}^{(L)} \left[\tilde{V}_{(L)}^{(\alpha_1 \alpha_3)(\alpha_2 \alpha_4)} + \frac{n_L}{2} \tilde{V}_{(L)}^{(\alpha_1 \alpha_2)(\alpha_3 \alpha_4)} \right]. \quad (\text{13})$$

Again, this can be derived through Wick's contractions. It is important to note that n_L is constant and thus $\epsilon \mathcal{H}_1[\mathbf{z}]$ can consistently be treated perturbatively.⁵

4. Interlude: examples

The recursive relations obtained above can be evaluated numerically (Lee et al., 2018) [or sometimes analytically for rectified linear unit (ReLU) activation (Cho and Saul, 2009)], which is a perfectly adequate approach: at the leading order it involves four-dimensional Gaussian integrals at most. Here, continuing the theme of wearing out Wick's contractions, we develop an alternative analytic method that works for any polynomial activations (Liao and Poggio, 2017), providing another perfectly cromulent approach.

For a general polynomial activation of degree p , $\sigma(z) = \sum_{k=0}^p a_k z^k$, the nontrivial term in Equation (R1) can be expanded as

$$\langle \sigma(\tilde{z}_{\alpha_1}) \sigma(\tilde{z}_{\alpha_2}) \rangle_{\tilde{K}^{(\ell)}} = \sum_{k_1, k_2=0}^p a_{k_1} a_{k_2} \left\langle (\tilde{z}_{\alpha_1})^{k_1} (\tilde{z}_{\alpha_2})^{k_2} \right\rangle_{\tilde{K}^{(\ell)}}. \quad (\text{14})$$

Each term can then be evaluated by Wick's contractions and the same goes for all the terms in Equations (R2) and (R3).⁶ Below and in Appendix C, we illustrate this procedure with simple examples.

4.1. Deep linear networks

When the activation function is linear, $\sigma(z) = z$, multilayer perceptrons are called deep linear networks (Saxe et al., 2013). Setting $C_b^{(\ell)} = 0$ and $C_W^{(\ell)} = 1$ for simplicity, our recursion relations reduce to $\tilde{K}_{\alpha_1 \alpha_2}^{(\ell+1)} = \tilde{K}_{\alpha_1 \alpha_2}^{(\ell)}$,

$$\tilde{V}_{(\alpha_1 \alpha_2)(\alpha_3 \alpha_4)}^{(\ell+1)} = \left[\tilde{K}_{\alpha_1 \alpha_3}^{(\ell)} \tilde{K}_{\alpha_2 \alpha_4}^{(\ell)} + \tilde{K}_{\alpha_1 \alpha_4}^{(\ell)} \tilde{K}_{\alpha_2 \alpha_3}^{(\ell)} + \left(\frac{n_\ell}{n_{\ell-1}} \right) \tilde{V}_{(\alpha_1 \alpha_2)(\alpha_3 \alpha_4)}^{(\ell)} \right],$$

5. If n_L were of order $n \gg 1$, the potential \mathcal{H} would become a large- n vector model, for which we would have to sum the infinite series of bubble diagrams (Moshe and Zinn-Justin, 2003).

6. The same approach could be adopted for an analytic function but it would in general be difficult to sum the resulting infinite series in a closed form. It could nonetheless be useful in, for example, proving convergence properties.

and $\tilde{S}_{\alpha_1\alpha_2}^{(\ell+1)} = \left(\frac{n_\ell}{n_{\ell-1}}\right) \tilde{S}_{\alpha_1\alpha_2}^{(\ell)}$. Solving them yields the layer-independent core kernel and zero self-energy

$$\tilde{K}_{\alpha_1\alpha_2}^{(\ell)} = \tilde{K}_{\alpha_1\alpha_2}^{(1)} = \frac{\mathbf{x}_{\alpha_1} \cdot \mathbf{x}_{\alpha_2}}{n_0} \quad \text{and} \quad \tilde{S}_{\alpha_1\alpha_2}^{(\ell)} = 0 \quad (15)$$

and the linearly layer-dependent four-point vertex

$$\frac{1}{n_{\ell-1}} \tilde{V}_{(\alpha_1\alpha_2)(\alpha_3\alpha_4)}^{(\ell)} = \left(\sum_{\ell'=1}^{\ell-1} \frac{1}{n_{\ell'}} \right) \left[\tilde{K}_{\alpha_1\alpha_3}^{(1)} \tilde{K}_{\alpha_2\alpha_4}^{(1)} + \tilde{K}_{\alpha_1\alpha_4}^{(1)} \tilde{K}_{\alpha_2\alpha_3}^{(1)} \right]. \quad (16)$$

It succinctly reproduces the result that can be obtained through planar diagrams in this special setup (Dyer and Gur-Ari, 2019). Quadratic activation (Li et al., 2018) is worked out in Appendix C.1.

4.2. ReLU with single input

The recursion relations simplify drastically for the case of a single input, $N_D = 1$, as worked out in detail in Appendix C.2. For instance, for ReLU activation with $C_b^{(\ell)} = 0$ and $C_W^{(\ell)} = 2$, we obtain the layer-independent core kernel, zero self-energy, and the four-point vertex

$$\frac{1}{n_{\ell-1}} \tilde{V}_{(\alpha\alpha)(\alpha\alpha)}^{(\ell)} = 5 \left(\sum_{\ell'=1}^{\ell-1} \frac{1}{n_{\ell'}} \right) \left(\tilde{K}_{\alpha\alpha}^{(1)} \right)^2. \quad (17)$$

Interestingly, as for deep linear networks, the factor $\sum_{\ell'} (1/n_{\ell'})$ appears again. This factor has also been found by Hanin and Rolnick (2018), which provides guidance for network architectural design through its minimization. We generalize this factor for monomial activations in Appendix C.2.1

4.3. Experimental verification: output distributions for a single input

Here we put our theory to the test. For concreteness, take a single black-white image of handwritten digits with 28-by-28 pixels (i.e. $n_0 = 784$) from the MNIST dataset (LeCun et al., 1998) without preprocessing, set depth $L = 3$, bias variance $C_b^{(\ell)} = 0$, weight variance $C_W^{(\ell)} = C_W$, and widths $(n_0, n_1, n_2, n_3) = (784, n, 2n, 1)$, and use activations $\sigma(z) = z$ (linear) with $C_W = 1$ and $\max(0, z)$ (ReLU) with $C_W = 2$. In Figure 1, for each width-parameter n of the hidden layers we record the prior distribution of outputs over 10^6 instances of Gaussian weights and compare it with the theoretical prediction – obtained by cranking the knob from the initial condition (R0) through the recursion relations (R1-R3) to the distribution (D0-D2). The prior distribution becomes increasingly non-Gaussian as networks narrow and the deviation from the Gaussian-process prior is correctly captured by our theory. Higher-order perturbative calculations are expected to systematically improve the quality – and extend the range – of the agreement. Additional experiments are performed in Appendix C.3, which further corroborates our theory.

5. Bayesian inference

Let us take off from the terminal point of Section 3: we have obtained the recursive equations (R0-R3) for the Gaussian-process kernel and the leading finite-width corrections and codified them in the weakly non-Gaussian prior distributions $p[\mathbf{z}]$ (D0-D2) of outputs

$$\mathbf{z} \equiv \left\{ z_{i;\alpha} \equiv z_i^{(L)}(\mathbf{x}_\alpha) \right\}_{i=1,\dots,n_L; \alpha=1,\dots,N_D},$$

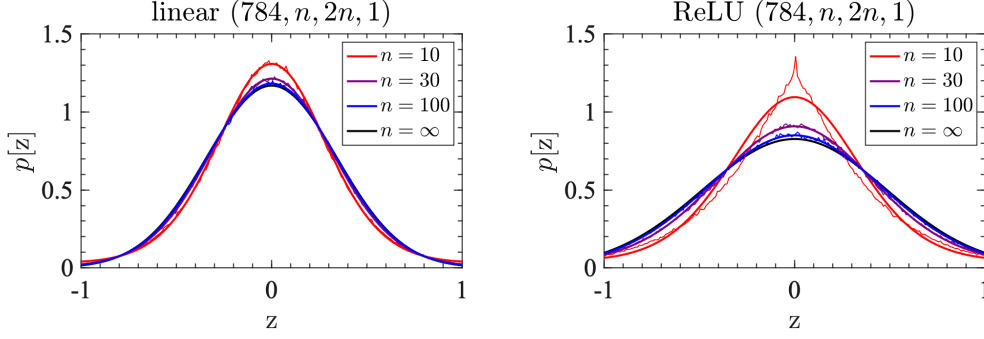


Figure 1: Comparison between theory and experiments for prior distributions of outputs for a single input. The agreement between our theoretical predictions (smooth thick lines) and experimental data (rugged thin lines) is superb, correctly capturing the initial deviations from Gaussian processes at $n = \infty$ (black), all the way down to $n \sim 10$ for linear activation and to $n \sim 30$ for ReLU activation.

dictated by the potential $\mathcal{H}[\mathbf{z}] = \mathcal{H}_0[\mathbf{z}] + \epsilon \mathcal{H}_1[\mathbf{z}] + O(\epsilon^2)$ with $\epsilon \equiv \frac{1}{n_{L-1}} \ll 1$. Examples in Section 4 illustrate that finite-width corrections stay perturbative typically when $\frac{\text{depth}}{\text{width}} \ll 1$. Let us now divide N_D inputs into N_R training and N_E test inputs as

$$\{\mathbf{x}_\alpha\}_{\alpha=1,\dots,N_D} = \{(\mathbf{x}_R)_{\bar{\beta}}\}_{\bar{\beta}=1,\dots,N_R} \cup \{(\mathbf{x}_E)_{\dot{\gamma}}\}_{\dot{\gamma}=1,\dots,N_E}, \quad (18)$$

and the training inputs come with target outputs

$$\{(\mathbf{y}_R)_{\bar{\beta}}\}_{\bar{\beta}=1,\dots,N_R} = \{(y_R)_{i;\bar{\beta}}\}_{\bar{\beta}=1,\dots,N_R; i=1,\dots,n_L}. \quad (19)$$

We shall develop a procedure to infer outputs for test inputs à la Bayes, perturbatively extending the textbook by [Williams and Rasmussen \(2006\)](#). For field theorists, our calculation is just a background-field calculation ([Weinberg, 1996](#)) in disguise.

Taking the liberty of notations, we let the number of input-data arguments dictate the summation over sample indices α inside the potential \mathcal{H} , and denote the joint probabilities

$$p[\mathbf{z}_R] = \frac{e^{-\mathcal{H}[\mathbf{z}_R]}}{\int d\mathbf{z}'_R e^{-\mathcal{H}[\mathbf{z}'_R]}} \quad \text{and} \quad p[\mathbf{z}_R, \mathbf{z}_E] = \frac{e^{-\mathcal{H}[\mathbf{z}_R, \mathbf{z}_E]}}{\int d\mathbf{z}'_R d\mathbf{z}'_E e^{-\mathcal{H}[\mathbf{z}'_R, \mathbf{z}'_E]}}. \quad (20)$$

Given the training targets \mathbf{y}_R , the posterior distribution of test outputs are given by Bayes' rule:

$$p[\mathbf{z}_E | \mathbf{y}_R] = \frac{p[\mathbf{y}_R, \mathbf{z}_E]}{p[\mathbf{y}_R]} = \left(\frac{\int d\mathbf{z}'_R e^{-\mathcal{H}[\mathbf{z}'_R]}}{\int d\mathbf{z}'_R d\mathbf{z}'_E e^{-\mathcal{H}[\mathbf{z}'_R, \mathbf{z}'_E]}} \right) e^{-(\mathcal{H}[\mathbf{y}_R, \mathbf{z}_E] - \mathcal{H}[\mathbf{y}_R])}. \quad (\text{Bayes})$$

The leading Gaussian-process contributions can be segregated out through the textbook manipulation ([Williams and Rasmussen, 2006](#)) [c.f. Appendix D]: denoting the full Gaussian-process kernel in the last layer as

$$K_{i_1 i_2; \alpha_1 \alpha_2} = \delta_{i_1 i_2} \begin{pmatrix} \left(\tilde{K}_{RR} \right)_{\bar{\beta}_1 \bar{\beta}_2} & \left(\tilde{K}_{RE} \right)_{\bar{\beta}_1 \dot{\gamma}_2} \\ \left(\tilde{K}_{ER} \right)_{\dot{\gamma}_1 \bar{\beta}_2} & \left(\tilde{K}_{EE} \right)_{\dot{\gamma}_1 \dot{\gamma}_2} \end{pmatrix} \quad (21)$$

and the Gaussian-process posterior mean prediction as

$$(\mathbf{y}_E^{\text{GP}})_{i;\dot{\gamma}} \equiv \sum_{\dot{\beta}} \left[\tilde{K}_{\text{ER}} \tilde{K}_{\text{RR}}^{-1} \right]_{\dot{\gamma}}^{\dot{\beta}} (\mathbf{y}_R)_{i;\dot{\beta}}, \quad (\text{GPM})$$

and defining a fluctuation $(\mathbf{z}_E)_{i;\dot{\gamma}} \equiv (\mathbf{y}_E^{\text{GP}})_{i;\dot{\gamma}} + (\delta \mathbf{z}_E)_{i;\dot{\gamma}}$ and a matrix $\tilde{K}_\Delta \equiv \tilde{K}_{\text{EE}} - \tilde{K}_{\text{ER}} \tilde{K}_{\text{RR}}^{-1} \tilde{K}_{\text{RE}}$,

$$\mathcal{H}_0[\mathbf{y}_R, \mathbf{z}_E] - \mathcal{H}_0[\mathbf{y}_R] = \frac{1}{2} \sum_i \sum_{\dot{\gamma}_1, \dot{\gamma}_2} (\delta \mathbf{z}_E)_{i;\dot{\gamma}_1} \left(\tilde{K}_\Delta^{-1} \right)^{\dot{\gamma}_1 \dot{\gamma}_2} (\delta \mathbf{z}_E)_{i;\dot{\gamma}_2}. \quad (\text{GPD})$$

For any function \mathcal{F} , its expectation over the Bayesian posterior (**Bayes**) then turns into

$$\int d\mathbf{z}_E \mathcal{F}[\mathbf{z}_E] p[\mathbf{z}_E | \mathbf{y}_R] = \tilde{\mathcal{N}} \left\langle e^{-\epsilon \mathcal{H}_1[\mathbf{y}_R, \mathbf{y}_E^{\text{GP}} + \delta \mathbf{z}_E]} \mathcal{F}[\mathbf{y}_E^{\text{GP}} + \delta \mathbf{z}_E] \right\rangle_{K_\Delta} \quad (22)$$

where the deviation kernel $\left\langle (\delta \mathbf{z}_E)_{i_1; \dot{\gamma}_1} (\delta \mathbf{z}_E)_{i_2; \dot{\gamma}_2} \right\rangle_{K_\Delta} \equiv \delta_{i_1 i_2} \left(\tilde{K}_\Delta \right)_{\dot{\gamma}_1 \dot{\gamma}_2}$ and the normalization factor

$$\tilde{\mathcal{N}} = \left[\left\langle e^{-\epsilon \mathcal{H}_1[\mathbf{y}_R, \mathbf{y}_E^{\text{GP}} + \delta \mathbf{z}_E]} \right\rangle_{K_\Delta} \right]^{-1} = 1 + O(\epsilon). \quad (23)$$

In particular the mean posterior output is given by

$$\begin{aligned} (\bar{\mathbf{y}}_E)_{i;\dot{\gamma}} &\equiv \int d\mathbf{z}_E (\mathbf{z}_E)_{i;\dot{\gamma}} p[\mathbf{z}_E | \mathbf{y}_R] = (\mathbf{y}_E^{\text{GP}})_{i;\dot{\gamma}} + \tilde{\mathcal{N}} \left\langle (\delta \mathbf{z}_E)_{i;\dot{\gamma}} e^{-\epsilon \mathcal{H}_1[\mathbf{y}_R, \mathbf{y}_E^{\text{GP}} + \delta \mathbf{z}_E]} \right\rangle_{K_\Delta} \\ &= (\mathbf{y}_E^{\text{GP}})_{i;\dot{\gamma}} - \epsilon \left\langle (\delta \mathbf{z}_E)_{i;\dot{\gamma}} \mathcal{H}_1[\mathbf{y}_R, \mathbf{y}_E^{\text{GP}} + \delta \mathbf{z}_E] \right\rangle_{K_\Delta} + O(\epsilon^2). \end{aligned} \quad (24)$$

Stringing together $\bar{\phi}_{i;\alpha} \equiv [(\mathbf{y}_R)_{i;\dot{\beta}}, (\mathbf{y}_E^{\text{GP}})_{i;\dot{\gamma}}]$, recalling Equation (D2) for \mathcal{H}_1 , and using Wick's contractions for one last time, the mean prediction becomes

$$\begin{aligned} &(\mathbf{y}_E^{\text{GP}})_{i;\dot{\gamma}} \quad (\text{NGPM}) \\ &+ \epsilon \sum_{\alpha_1, \dot{\gamma}_1} \left(\tilde{K}_\Delta \right)_{\dot{\gamma} \dot{\gamma}_1} \bar{\phi}_{i;\alpha_1} \left[\tilde{S}^{\dot{\gamma}_1 \alpha_1} - \sum_{\alpha_2, \alpha_3} \tilde{V}^{(\dot{\gamma}_1 \alpha_2)(\alpha_1 \alpha_3)} \tilde{K}_{\alpha_2 \alpha_3} + \sum_{\dot{\gamma}_2, \dot{\gamma}_3} \tilde{V}^{(\dot{\gamma}_1 \dot{\gamma}_2)(\alpha_1 \dot{\gamma}_3)} \left(\tilde{K}_\Delta \right)_{\dot{\gamma}_2 \dot{\gamma}_3} \right. \\ &\left. + \frac{n_L}{2} \sum_{\dot{\gamma}_2, \dot{\gamma}_3} \tilde{V}^{(\alpha_1 \dot{\gamma}_1)(\dot{\gamma}_2 \dot{\gamma}_3)} \left(\tilde{K}_\Delta \right)_{\dot{\gamma}_2 \dot{\gamma}_3} + \sum_{\alpha_2, \alpha_3} \tilde{V}^{(\dot{\gamma}_1 \alpha_1)(\alpha_2 \alpha_3)} \left(-\frac{n_L}{2} \tilde{K}_{\alpha_2 \alpha_3} + \frac{1}{2} \sum_j \bar{\phi}_{j;\alpha_2} \bar{\phi}_{j;\alpha_3} \right) \right]. \end{aligned}$$

With additional manipulations, this expression can be simplified into the actionable form that is amenable to use in practice [c.f. Equations (NGPM') and (NGPM'') in Appendix D]. It turns out that for deep linear networks the leading finite-width correction vanishes, and the first correction is likely to show up at higher order in $1/n$ asymptotic expansion, which is not carried out in this paper. Here we instead use the $L = 2$ multilayer perceptron with the quadratic activation $\sigma(z) = z^2$, zero bias variance $C_b^{(\ell)} = 0$, and weight variance $C_W^{(\ell)} = 1/3$ for illustration, plugging Equations (S30,S31,S32) into Equations (NGPM') and (NGPM'') and varying $\epsilon \equiv \frac{1}{n_{L-1}} = \frac{1}{n_1}$. Results in Figure 2 indicate the regularization effects of finite widths when the number of training samples, N_R , is small, resulting in peak performance at finite widths. This is in line with expectations that finite widths ameliorate overfitting and that non-Gaussian priors increase the expressivity of neural functions, but additional large-scale extensive experiments would be desirable in the future.

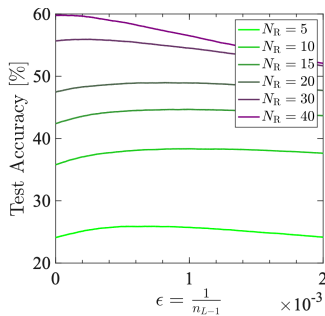


Figure 2: Test accuracy for $N_E = 10000$ MNIST test data as a function of the inverse width $\epsilon = 1/n_{L-1}$ of the hidden layer with quadratic activation. For each number N_R of subsampled training data, the result is averaged over 10 distinct choices of such subsamplings. For small numbers of training data, finite widths result in regularization effects, improving the test accuracy.

6. Dreams

In this paper, we have developed the perturbative formalism that captures the flow of preactivation distributions from lower to higher layers. The resemblance between our recursive equations and renormalization-group flow equations in high-energy and statistical physics is highly appealing. It would be exciting to investigate the structure of fixed points away from the Gaussian asymptopia (Schoenholz et al., 2016) and fully realize the dream articulated by Mehta and Schwab (2014) – the audacious hypothesis that neural networks wash away microscopic irrelevancies and extract relevant features – beyond their limited example of a mapping between two antiquated techniques.

In addition we have developed the perturbative Bayesian inference scheme universally applicable whenever prior distributions are weakly non-Gaussian, and have applied it to the specific cases of neural networks at finite widths. In light of possible finite-width regularization effects, it would be prudent to revisit the empirical comparison between SGD optimization and Bayesian inference at finite widths (Lee et al., 2018; Novak et al., 2019), especially for convolutional neural networks.

Finally, given surging interests in SGD dynamics within the large-width regime (Jacot et al., 2018; Chizat et al., 2018; Lee et al., 2019; Cohen et al., 2019; Dyer and Gur-Ari, 2019), it would be natural to adapt our formalism for investigating corrections to neural tangent kernels, and even aspire to capture a transition from lazy-learning to feature-learning regimes.

Acknowledgments

The author thanks Yasaman Bahri for the discussion that seeded the idea for this project, Boris L. Hanin for persistently preaching about Gaussian processes, and David J. Schwab for permission to call his example limited with our friendship intact. The author also thanks Ethan S. Dyer, Mario Geiger, Guy Gur-Ari, Eric T. Mintun, Stephen H. Shenker, and Lexing Ying for substantially useful discussions, and Daniel A. Roberts for the quality control of all the jokes and more.

References

- Joseph M. Antognini. Finite size corrections for neural network Gaussian processes. *arXiv preprint arXiv:1908.10030*, 2019.
- Robert Brown. XXVII. A brief account of microscopical observations made in the months of June, July and August 1827, on the particles contained in the pollen of plants; and on the general existence of active molecules in organic and inorganic bodies. *The Philosophical Magazine*, 4 (21):161–173, 1828.
- Lenaïc Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. *arXiv preprint arXiv:1812.07956*, 2018.
- Youngmin Cho and Lawrence K. Saul. Kernel methods for deep learning. In *Advances in Neural Information Processing Systems*, pages 342–350, 2009.
- Omry Cohen, Or Malka, and Zohar Ringel. Learning curves for deep neural networks: a Gaussian field theory perspective. *arXiv preprint arXiv:1906.05301*, 2019.
- Ethan Dyer and Guy Gur-Ari. Asymptotics of wide networks from Feynman diagrams. *arXiv preprint arXiv:1909.11304*, 2019.
- Albert Einstein. Über die von der molekularkinetischen Theorie der Wärme geforderte Bewegung von in ruhenden Flüssigkeiten suspendierten Teilchen. *Annalen der Physik*, 322(8):549–560, 1905.
- Mario Geiger, Stefano Spigler, Arthur Jacot, and Matthieu Wyart. Disentangling feature and lazy learning in deep neural networks: an empirical study. *arXiv preprint arXiv:1906.08034*, 2019.
- Nigel Goldenfeld. *Lectures on phase transitions and the renormalization group*. CRC Press, 2018.
- Boris Hanin and David Rolnick. How to start training: the effect of initialization and architecture. In *Advances in Neural Information Processing Systems*, pages 571–581, 2018.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems*, pages 8571–8580, 2018.
- Leo P. Kadanoff. Scaling laws for Ising models near T_c . *Physics Physique Fizika*, 2(6):263–272, 1966.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S. Schoenholz, Jeffrey Pennington, and Jascha Sohl-Dickstein. Deep neural networks as Gaussian processes. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=B1EA-M-0Z>.

- Jaehoon Lee, Lechao Xiao, Samuel S. Schoenholz, Yasaman Bahri, Jascha Sohl-Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models under gradient descent. *arXiv preprint arXiv:1902.06720*, 2019.
- Yuanzhi Li, Tengyu Ma, and Hongyang Zhang. Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations. In *Conference On Learning Theory*, pages 2–47, 2018.
- Qianli Liao and Tomaso Poggio. Theory II: landscape of the empirical risk in deep learning. *arXiv preprint arXiv:1703.09833*, 2017.
- David J. C. MacKay. Probable networks and plausible predictions – a review of practical Bayesian methods for supervised neural networks. *Network: Computation in Neural Systems*, 6(3):469–505, 1995.
- Alexander G. de G. Matthews, Mark Rowland, Jiri Hron, Richard E. Turner, and Zoubin Ghahramani. Gaussian process behaviour in wide deep neural networks. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=H1-nGgWC->.
- Pankaj Mehta and David J. Schwab. An exact mapping between the variational renormalization group and deep learning. *arXiv preprint arXiv:1410.3831*, 2014.
- Moshe Moshe and Jean Zinn-Justin. Quantum field theory in the large N limit: a review. *Physics Reports*, 385(3-6):69–228, 2003.
- Radford M. Neal. Priors for infinite networks. In *Bayesian Learning for Neural Networks*, pages 29–53. Springer, 1996.
- Roman Novak, Lechao Xiao, Yasaman Bahri, Jaehoon Lee, Greg Yang, Jiri Hron, Daniel A. Abolafia, Jeffrey Pennington, and Jascha Sohl-Dickstein. Bayesian deep convolutional networks with many channels are Gaussian processes. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=B1g30j0qF7>.
- Giorgio Parisi and Francesco Zamponi. Mean-field theory of hard sphere glasses and jamming. *Reviews of Modern Physics*, 82(1):789–845, 2010.
- Herbert Robbins and Sutton Monro. A stochastic approximation method. *The Annals of Mathematical Statistics*, 22(3):400–407, 1951.
- Andrew M. Saxe, James L. McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013.
- Samuel S. Schoenholz, Justin Gilmer, Surya Ganguli, and Jascha Sohl-Dickstein. Deep information propagation. In *International Conference on Learning Representations*, 2016. URL <https://openreview.net/forum?id=H1W1UN9gg>.
- Steven Weinberg. *The quantum theory of fields: Volume I: Foundations*. Cambridge University Press, 1995.

Steven Weinberg. *The quantum theory of fields: Volume II: Modern Applications*. Cambridge University Press, 1996.

Gian-Carlo Wick. The evaluation of the collision matrix. *Physical Review*, 80(2):268–272, 1950.

Christopher K. I. Williams. Computing with infinite networks. In *Advances in Neural Information Processing Systems*, pages 295–301, 1997.

Christopher K. I. Williams and Carl Edward Rasmussen. *Gaussian processes for machine learning*. MIT Press, 2006.

Kenneth G. Wilson. Renormalization group and critical phenomena. I. Renormalization group and the Kadanoff scaling picture. *Physical Review B*, 4(9):3174–3183, 1971.

Anthony Zee. *Quantum field theory in a nutshell*. Princeton University Press, 2010.

Appendix A. Wick's tricks

Here is all you need to know in order to follow the calculations in the paper. In the main text, Wick's contractions are used both for trivially integrating out biases and weights as straightforward applications of Appendix A.1 and for nontrivially integrating out preactivations, with concepts of cumulants reviewed in Appendix A.2 and A.3, culminating in the hack derived in Appendix A.4. The random variables are generically indexed by $\mu = 1, \dots, N$ throughout this Appendix: when applying formulae for biases, $\mu = i$; for weights $\mu = (i, j)$; for full preactivations $\mu = (i, \alpha)$; for single-neuron preactivations $\mu = \alpha$.

A.1. Wick's contractions

For Gaussian-distributed variables $\mathbf{z} = \{z_\mu\}_{\mu=1, \dots, N}$ with a kernel $K_{\mu\mu'}$, moments

$$\langle z_{\mu_1} z_{\mu_2} \cdots z_{\mu_m} \rangle_K \equiv \frac{\int d\mathbf{z} e^{-\mathcal{H}_0[\mathbf{z}]} z_{\mu_1} z_{\mu_2} \cdots z_{\mu_m}}{\int d\mathbf{z} e^{-\mathcal{H}_0[\mathbf{z}]}} \quad \text{with} \quad \mathcal{H}_0[\mathbf{z}] \equiv \frac{1}{2} \sum_{\mu, \mu'=1}^N z_\mu (K^{-1})^{\mu\mu'} z_{\mu'}. \quad (\text{S1})$$

For any odd m such moments identically vanish. For even m , Isserlis-Wick's theorem states that

$$\langle z_{\mu_1} z_{\mu_2} \cdots z_{\mu_m} \rangle_K = \sum_{\text{all pairing}} K_{\mu_{k_1} \mu_{k_2}} \cdots K_{\mu_{k_{m-1}} \mu_{k_m}} \quad (\text{S2})$$

where the sum is over all the possible pairings of m variables, $(k_1, k_2), \dots, (k_{m-1}, k_m)$. In general, there are $(m-1)!! = (m-1) \cdot (m-3) \cdots 1$ such pairings. For a proof, see for example Zee (2010). In order to understand and use the theorem, it is instructive to look at a few examples:

$$\langle z_{\mu_1} z_{\mu_2} \rangle_K = K_{\mu_1 \mu_2}; \quad (\text{S3})$$

$$\langle z_{\mu_1} z_{\mu_2} z_{\mu_3} z_{\mu_4} \rangle_K = K_{\mu_1 \mu_2} K_{\mu_3 \mu_4} + K_{\mu_1 \mu_3} K_{\mu_2 \mu_4} + K_{\mu_1 \mu_4} K_{\mu_2 \mu_3}; \quad (\text{S4})$$

and

$$\begin{aligned} & \langle z_{\mu_1} z_{\mu_2} z_{\mu_3} z_{\mu_4} z_{\mu_5} z_{\mu_6} \rangle_K \quad (\text{S5}) \\ &= K_{\mu_1 \mu_2} K_{\mu_3 \mu_4} K_{\mu_5 \mu_6} + K_{\mu_1 \mu_3} K_{\mu_2 \mu_4} K_{\mu_5 \mu_6} + K_{\mu_1 \mu_4} K_{\mu_2 \mu_3} K_{\mu_5 \mu_6} \\ &+ K_{\mu_1 \mu_2} K_{\mu_3 \mu_5} K_{\mu_4 \mu_6} + K_{\mu_1 \mu_3} K_{\mu_2 \mu_5} K_{\mu_4 \mu_6} + K_{\mu_1 \mu_5} K_{\mu_2 \mu_3} K_{\mu_4 \mu_6} \\ &+ K_{\mu_1 \mu_2} K_{\mu_3 \mu_4} K_{\mu_5 \mu_6} + K_{\mu_1 \mu_5} K_{\mu_2 \mu_4} K_{\mu_3 \mu_6} + K_{\mu_1 \mu_4} K_{\mu_2 \mu_5} K_{\mu_3 \mu_6} \\ &+ K_{\mu_1 \mu_5} K_{\mu_3 \mu_4} K_{\mu_2 \mu_6} + K_{\mu_1 \mu_3} K_{\mu_5 \mu_4} K_{\mu_2 \mu_6} + K_{\mu_1 \mu_4} K_{\mu_5 \mu_3} K_{\mu_2 \mu_6} \\ &+ K_{\mu_5 \mu_2} K_{\mu_3 \mu_4} K_{\mu_1 \mu_6} + K_{\mu_5 \mu_3} K_{\mu_2 \mu_4} K_{\mu_1 \mu_6} + K_{\mu_5 \mu_4} K_{\mu_2 \mu_3} K_{\mu_1 \mu_6}. \end{aligned}$$

A.2. Connected correlations

Given general (not necessarily Gaussian) random variables, connected correlation functions are defined inductively through

$$\begin{aligned} & \mathbb{E} [z_{\mu_1} z_{\mu_2} \cdots z_{\mu_m}] \quad (\text{S6}) \\ & \equiv \mathbb{E} [z_{\mu_1} z_{\mu_2} \cdots z_{\mu_m}] \Big|_{\text{connected}} \\ & + \sum_{\text{all subdivisions}} \mathbb{E} \left[z_{\mu_{k_1}^{[1]}} \cdots z_{\mu_{k_{l_1}^{[1]}}^{[1]}} \right] \Big|_{\text{connected}} \cdots \mathbb{E} \left[z_{\mu_{k_1^{[s]}}} \cdots z_{\mu_{k_{l_s^{[s]}}}^{[s]}} \right] \Big|_{\text{connected}} \end{aligned}$$

where the sum is over all the possible subdivisions of m variables into $s > 1$ clusters of sizes (ν_1, \dots, ν_s) as $(k_1^{[1]}, \dots, k_{\nu_1}^{[1]}), \dots, (k_1^{[s]}, \dots, k_{\nu_s}^{[s]})$. In order to understand the definition, it is again instructive to look at a few examples. Assuming that all the odd moments vanish,

$$\mathbb{E}[z_{\mu_1} z_{\mu_2}] = \mathbb{E}[z_{\mu_1} z_{\mu_2}] \Big|_{\text{connected}} \quad \text{and} \quad (\text{S7})$$

$$\begin{aligned} \mathbb{E}[z_{\mu_1} z_{\mu_2} z_{\mu_3} z_{\mu_4}] &= \mathbb{E}[z_{\mu_1} z_{\mu_2} z_{\mu_3} z_{\mu_4}] \Big|_{\text{connected}} \\ &+ \mathbb{E}[z_{\mu_1} z_{\mu_2}] \Big|_{\text{connected}} \mathbb{E}[z_{\mu_3} z_{\mu_4}] \Big|_{\text{connected}} \\ &+ \mathbb{E}[z_{\mu_1} z_{\mu_3}] \Big|_{\text{connected}} \mathbb{E}[z_{\mu_2} z_{\mu_4}] \Big|_{\text{connected}} \\ &+ \mathbb{E}[z_{\mu_1} z_{\mu_4}] \Big|_{\text{connected}} \mathbb{E}[z_{\mu_2} z_{\mu_3}] \Big|_{\text{connected}} . \end{aligned} \quad (\text{S8})$$

Rearranging them in particular yields

$$\begin{aligned} &\mathbb{E}[z_{\mu_1} z_{\mu_2} z_{\mu_3} z_{\mu_4}] \Big|_{\text{connected}} \\ &= \mathbb{E}[z_{\mu_1} z_{\mu_2} z_{\mu_3} z_{\mu_4}] \\ &\quad - \mathbb{E}[z_{\mu_1} z_{\mu_2}] \mathbb{E}[z_{\mu_3} z_{\mu_4}] - \mathbb{E}[z_{\mu_1} z_{\mu_3}] \mathbb{E}[z_{\mu_2} z_{\mu_4}] - \mathbb{E}[z_{\mu_1} z_{\mu_4}] \mathbb{E}[z_{\mu_2} z_{\mu_3}] . \end{aligned} \quad (\text{S9})$$

If these examples do not suffice, here is yet another example to chew on:

$$\begin{aligned} \mathbb{E}[z_{\mu_1} z_{\mu_2} z_{\mu_3} z_{\mu_4} z_{\mu_5} z_{\mu_6}] &= \mathbb{E}[z_{\mu_1} z_{\mu_2} z_{\mu_3} z_{\mu_4} z_{\mu_5} z_{\mu_6}] \Big|_{\text{connected}} \\ &+ \mathbb{E}[z_{\mu_1} z_{\mu_2}] \Big|_{\text{connected}} \mathbb{E}[z_{\mu_3} z_{\mu_4}] \Big|_{\text{connected}} \mathbb{E}[z_{\mu_5} z_{\mu_6}] \Big|_{\text{connected}} \\ &+ [14 \text{ other } (2, 2, 2) \text{ subdivisions}] \\ &+ \mathbb{E}[z_{\mu_1} z_{\mu_2} z_{\mu_3} z_{\mu_4}] \Big|_{\text{connected}} \mathbb{E}[z_{\mu_5} z_{\mu_6}] \Big|_{\text{connected}} \\ &+ [14 \text{ other } (4, 2) \text{ subdivisions}] \end{aligned} \quad (\text{S10})$$

and hence

$$\begin{aligned} &\mathbb{E}[z_{\mu_1} z_{\mu_2} z_{\mu_3} z_{\mu_4} z_{\mu_5} z_{\mu_6}] \Big|_{\text{connected}} \\ &= \mathbb{E}[z_{\mu_1} z_{\mu_2} z_{\mu_3} z_{\mu_4} z_{\mu_5} z_{\mu_6}] \\ &\quad - \{\mathbb{E}[z_{\mu_1} z_{\mu_2} z_{\mu_3} z_{\mu_4}] \mathbb{E}[z_{\mu_5} z_{\mu_6}] + [14 \text{ other } (4, 2) \text{ subdivisions}]\} \\ &\quad + 2 \{\mathbb{E}[z_{\mu_1} z_{\mu_2}] \mathbb{E}[z_{\mu_3} z_{\mu_4}] \mathbb{E}[z_{\mu_5} z_{\mu_6}] + [14 \text{ other } (2, 2, 2) \text{ subdivisions}]\} . \end{aligned} \quad (\text{S11})$$

We emphasize that these are just renderings of the definition (S6). The power of this definition will be illustrated in the next two subsections.

A.3. Hierarchical clustering

We often encounter situations with the hierarchy

$$\mathbb{E}[z_{\mu_1} \cdots z_{\mu_m}] \Big|_{\text{connected}} = O(\epsilon^{\frac{m-2}{2}}) \quad (\text{S12})$$

where $\epsilon \ll 1$ is a small perturbative parameter and here again odd moments are assumed to vanish. Often comes with the hierarchical structure is the asymptotic limit $\epsilon \rightarrow 0$ where

$$\mathbb{E}[z_{\mu_1} z_{\mu_2}] = K_{\mu_1 \mu_2} + \epsilon S_{\mu_1 \mu_2} + O(\epsilon^2) \quad (\text{S13})$$

with the Gaussian kernel $K_{\mu_1\mu_2}$ at zero ϵ and the leading self-energy correction $S_{\mu_1\mu_2}$. Let us also denote the leading four-point vertex

$$\mathbb{E} [z_{\mu_1} z_{\mu_2} z_{\mu_3} z_{\mu_4}] \Big|_{\text{connected}} = \epsilon V_{\mu_1\mu_2\mu_3\mu_4} + O(\epsilon^2). \quad (\text{S14})$$

For instance this hierarchy holds for weakly-coupled field theories – from which we are importing names such as self-energy, vertex, and metric – and, in this paper, such hierarchical structure is inductively shown to hold for prior preactivations $\mathbf{z}^{(\ell)}$ with $\epsilon = \frac{1}{n_{\ell-1}}$ in the regime $n_1, \dots, n_{L-1} \sim n \gg 1$. Note that, by definition, $K_{\mu_1\mu_2}$ and $S_{\mu_1\mu_2}$ are symmetric under $\mu_1 \leftrightarrow \mu_2$ and $V_{\mu_1\mu_2\mu_3\mu_4}$ is symmetric under permutations of $(\mu_1, \mu_2, \mu_3, \mu_4)$.⁷

A.4. Combinatorial hack

So far we have reviewed the standard technology of Wick’s contractions, connected correlation functions, and all that. Our objective now is to develop a method to evaluate $\mathbb{E} [z_{\mu_1} \cdots z_{\mu_m}]$ for random variables obeying the hierarchical-clustering property (S12),⁸ which is the inductive hypothesis made in the main text; by extension, the resulting method (HACK’) lets us perturbatively evaluate $\mathbb{E} \{\mathcal{F}[\mathbf{z}]\}$ for any function \mathcal{F} that can be obtained as a limit of a sequence of analytic functions.

With the review of connected correlation functions passed us, first note that

$$\begin{aligned} & \mathbb{E} [z_{\mu_1} \cdots z_{\mu_m}] && (\text{CLUSTER}) \\ &= \mathbb{E} [z_{\mu_1} z_{\mu_2}] \cdots \mathbb{E} [z_{\mu_{m-1}} z_{\mu_m}] + \{[(m-1)!! - 1] \text{ other pairings}\} \\ & \quad + \mathbb{E} [z_{\mu_1} z_{\mu_2} z_{\mu_3} z_{\mu_4}] \Big|_{\text{connected}} \mathbb{E} [z_{\mu_5} z_{\mu_6}] \cdots \mathbb{E} [z_{\mu_{m-1}} z_{\mu_m}] \\ & \quad + \left\{ \left[\binom{m}{4} \times (m-5)!! - 1 \right] \text{ other } (4, 2, 2, \dots, 2) \text{ clusterings} \right\} + O(\epsilon^2) \\ &= K_{\mu_1\mu_2} \cdots K_{\mu_{m-1}\mu_m} + \{[(m-1)!! - 1] \text{ other pairings}\} \\ & \quad + \epsilon S_{\mu_1\mu_2} K_{\mu_3\mu_4} \cdots K_{\mu_{m-1}\mu_m} \\ & \quad + \left\{ \left[\binom{m}{2} \times (m-3)!! - 1 \right] \text{ other such clusterings} \right\} \\ & \quad + \epsilon V_{\mu_1\mu_2\mu_3\mu_4} K_{\mu_5\mu_6} \cdots K_{\mu_{m-1}\mu_m} \\ & \quad + \left\{ \left[\binom{m}{4} \times (m-5)!! - 1 \right] \text{ other } (4, 2, 2, \dots, 2) \text{ clusterings} \right\} + O(\epsilon^2) \\ &= \langle z_{\mu_1} \cdots z_{\mu_m} \rangle_K && (\text{CLUSTER}') \\ & \quad + \epsilon S_{\mu_1\mu_2} \langle z_{\mu_3} \cdots z_{\mu_m} \rangle_K + \left\{ \left[\binom{m}{2} - 1 \right] \text{ other self-energy contractions} \right\} \\ & \quad + \epsilon V_{\mu_1\mu_2\mu_3\mu_4} \langle z_{\mu_5} \cdots z_{\mu_m} \rangle_K + \left\{ \left[\binom{m}{4} - 1 \right] \text{ other vertex contractions} \right\} + O(\epsilon^2). \end{aligned}$$

where in the last equality Wick’s theorem was used backward.

7. In the main text the connected four-point preactivation correlation functions are symmetric under the permutations of four (sample, neuron) indices, $\{(i_1, \alpha_1), (i_2, \alpha_2), (i_3, \alpha_3), (i_4, \alpha_4)\}$.

8. More precisely, we shall inductively use only the weaker proposition that $\mathbb{E} [z_{\mu_1} \cdots z_{\mu_m}] \Big|_{\text{connected}} = O(\epsilon^2)$ for $m \geq 6$ along with Equations (S13) and (S14).

Below, let us use the inverse kernel $(K^{-1})^{\mu_1\mu_2}$ as a metric to raise indices:

$$S^{\mu_1\mu_2} \equiv \sum_{\mu'_1, \mu'_2} (K^{-1})^{\mu_1\mu'_1} (K^{-1})^{\mu_2\mu'_2} S_{\mu'_1\mu'_2} \quad \text{and} \quad (\text{S15})$$

$$V^{\mu_1\mu_2\mu_3\mu_4} \equiv \sum_{\mu'_1, \dots, \mu'_4} (K^{-1})^{\mu_1\mu'_1} \dots (K^{-1})^{\mu_4\mu'_4} V_{\mu'_1\mu'_2\mu'_3\mu'_4}. \quad (\text{S16})$$

Then, in order to simplify the second set of terms in Equation (CLUSTER') involving self-energy, note that

$$\begin{aligned} & \left\langle \mathbf{z}_{\mu_1} \cdots \mathbf{z}_{\mu_m} \left(\sum_{\mu'_1, \mu'_2} S^{\mu'_1\mu'_2} \mathbf{z}_{\mu'_1} \mathbf{z}_{\mu'_2} \right) \right\rangle_K \\ &= \sum_{\mu'_1, \mu'_2} S^{\mu'_1\mu'_2} \left[\left\langle \mathbf{z}_{\mu'_1} \mathbf{z}_{\mu'_2} \right\rangle_K \left\langle \mathbf{z}_{\mu_1} \cdots \mathbf{z}_{\mu_m} \right\rangle_K \right. \\ & \quad \left. + 2 \left\langle \mathbf{z}_{\mu_1} \mathbf{z}_{\mu'_1} \right\rangle_K \left\langle \mathbf{z}_{\mu_2} \mathbf{z}_{\mu'_2} \right\rangle_K \left\langle \mathbf{z}_{\mu_3} \cdots \mathbf{z}_{\mu_m} \right\rangle_K + \left\{ \left[\binom{m}{2} - 1 \right] \text{other } (\mu_1, \mu_2) \right\} \right] \\ &= \left(\sum_{\mu'_1, \mu'_2} S^{\mu'_1\mu'_2} K_{\mu'_1\mu'_2} \right) \left\langle \mathbf{z}_{\mu_1} \cdots \mathbf{z}_{\mu_m} \right\rangle_K \\ & \quad + 2S_{\mu_1\mu_2} \left\langle \mathbf{z}_{\mu_3} \cdots \mathbf{z}_{\mu_m} \right\rangle_K + \left\{ \left[\binom{m}{2} - 1 \right] \text{other } (\mu_1, \mu_2) \right\} \end{aligned}$$

where the symmetry $\mu_1 \leftrightarrow \mu_2$ of $S_{\mu_1\mu_2}$ was used. Hence, defining

$$\mathcal{O}_S[\mathbf{z}] \equiv \frac{1}{2} \sum_{\mu'_1, \mu'_2} S^{\mu'_1\mu'_2} \left(\mathbf{z}_{\mu'_1} \mathbf{z}_{\mu'_2} - K_{\mu'_1\mu'_2} \right), \quad (\text{OS}')$$

we obtain

$$\epsilon S_{\mu_1\mu_2} \left\langle \mathbf{z}_{\mu_3} \cdots \mathbf{z}_{\mu_m} \right\rangle_K + \left\{ \left[\binom{m}{2} - 1 \right] \text{other } (\mu_1, \mu_2) \right\} \quad (\text{S17})$$

$$= \epsilon \left\langle \mathbf{z}_{\mu_1} \cdots \mathbf{z}_{\mu_m} \mathcal{O}_S[\mathbf{z}] \right\rangle_K. \quad (\text{S18})$$

The similar algebraic exercise renders the other term in Equation (CLUSTER') to be

$$\epsilon V_{\mu_1\mu_2\mu_3\mu_4} \left\langle \mathbf{z}_{\mu_5} \cdots \mathbf{z}_{\mu_m} \right\rangle_K + \left\{ \left[\binom{m}{4} - 1 \right] \text{other } (\mu_1, \mu_2, \mu_3, \mu_4) \right\} \quad (\text{S19})$$

$$= \epsilon \left\langle \mathbf{z}_{\mu_1} \cdots \mathbf{z}_{\mu_m} \mathcal{O}_V[\mathbf{z}] \right\rangle_K \quad (\text{S20})$$

with

$$\mathcal{O}_V[\mathbf{z}] \equiv \frac{1}{24} \sum_{\mu'_1, \dots, \mu'_4} V^{\mu'_1\mu'_2\mu'_3\mu'_4} \left(\mathbf{z}_{\mu'_1} \mathbf{z}_{\mu'_2} \mathbf{z}_{\mu'_3} \mathbf{z}_{\mu'_4} - 6\mathbf{z}_{\mu'_1} \mathbf{z}_{\mu'_2} K_{\mu'_3\mu'_4} + 3K_{\mu'_1\mu'_2} K_{\mu'_3\mu'_4} \right). \quad (\text{OV}')$$

In summary, for any function $\mathcal{F}[\mathbf{z}]$ of random variables z_μ

$$\mathbb{E} \{ \mathcal{F}[\mathbf{z}] \} = \langle \mathcal{F}[\mathbf{z}] \rangle_K + \epsilon \langle \mathcal{F}[\mathbf{z}] \mathcal{O}_S[\mathbf{z}] \rangle_K + \epsilon \langle \mathcal{F}[\mathbf{z}] \mathcal{O}_V[\mathbf{z}] \rangle_K + O(\epsilon^2). \quad (\text{HACK'})$$

In order to get the expressions used in the main text at the ℓ -th layer, we need only to replace $\mu \rightarrow (i, \alpha)$, identify $\epsilon = \frac{1}{n_{\ell-1}}$, and use the inductive hypotheses **(KS)**

$$\mathbb{E} [z_{i_1; \alpha_1} z_{i_2; \alpha_2}] = \delta_{i_1 i_2} \left[\tilde{K}_{\alpha_1 \alpha_2}^{(\ell)} + \frac{1}{n_{\ell-1}} \tilde{S}_{\alpha_1 \alpha_2}^{(\ell)} + O\left(\frac{1}{n^2}\right) \right]$$

and **(V)**

$$\begin{aligned} & \mathbb{E} [z_{i_1; \alpha_1} z_{i_2; \alpha_2} z_{i_3; \alpha_3} z_{i_4; \alpha_4}] \Big|_{\text{connected}} \\ &= \frac{1}{n_{\ell-1}} \left[\delta_{i_1 i_2} \delta_{i_3 i_4} \tilde{V}_{(\alpha_1 \alpha_2)(\alpha_3 \alpha_4)}^{(\ell)} + \delta_{i_1 i_3} \delta_{i_2 i_4} \tilde{V}_{(\alpha_1 \alpha_3)(\alpha_2 \alpha_4)}^{(\ell)} \right. \\ & \quad \left. + \delta_{i_1 i_4} \delta_{i_2 i_3} \tilde{V}_{(\alpha_1 \alpha_4)(\alpha_2 \alpha_3)}^{(\ell)} \right] + O\left(\frac{1}{n^2}\right). \end{aligned}$$

The operators in Equations **(OS')** and **(OV')** then become

$$\begin{aligned} \mathcal{O}_S[\mathbf{z}] &= \frac{1}{2} \sum_{\alpha_1, \alpha_2} \tilde{S}_{(\ell)}^{\alpha_1 \alpha_2} \left[\left(\sum_{i=1}^{n_\ell} z_{i; \alpha_1} z_{i; \alpha_2} \right) - n_\ell \tilde{K}_{\alpha_1 \alpha_2}^{(\ell)} \right] \quad \text{and} \\ \mathcal{O}_V[\mathbf{z}] &= \frac{1}{8} \sum_{\alpha_1, \alpha_2, \alpha_3, \alpha_4} \tilde{V}_{(\ell)}^{(\alpha_1 \alpha_2)(\alpha_3 \alpha_4)} \\ & \quad \times \left[\left(\sum_{i=1}^{n_\ell} z_{i; \alpha_1} z_{i; \alpha_2} \right) \left(\sum_{j=1}^{n_\ell} z_{j; \alpha_3} z_{j; \alpha_4} \right) - 2n_\ell \left(\sum_{i=1}^{n_\ell} z_{i; \alpha_1} z_{i; \alpha_2} \right) \tilde{K}_{\alpha_3 \alpha_4}^{(\ell)} \right. \\ & \quad \left. - 4 \left(\sum_{i=1}^{n_\ell} z_{i; \alpha_1} z_{i; \alpha_3} \right) \tilde{K}_{\alpha_2 \alpha_4}^{(\ell)} + n_\ell^2 \tilde{K}_{\alpha_1 \alpha_2}^{(\ell)} \tilde{K}_{\alpha_3 \alpha_4}^{(\ell)} + 2n_\ell \tilde{K}_{\alpha_1 \alpha_3}^{(\ell)} \tilde{K}_{\alpha_2 \alpha_4}^{(\ell)} \right], \end{aligned}$$

i.e., the operators in Equations **(OS)** and **(OV)** in the main text.

Appendix B. Full condensed proof

In this Appendix, we provide a full inductive proof for one of the main claims in the paper, streamlined in the main text. Namely, we assume at the ℓ -th layer that Equations (KS) and (V) hold and that all the higher-point connected preactivation correlation functions are of order $O\left(\frac{1}{n^2}\right)$ – which are trivially true at $\ell = 1$ – and prove the same for the $(\ell + 1)$ -th layer. We assume the full mastery of Appendix A or, conversely, this section can be used to test the mastery of Wick’s tricks.

First, trivial Wick’s contractions yield

$$\begin{aligned}
 & G_{i_1 \dots i_{2m}; \alpha_1 \dots \alpha_{2m}}^{(\ell+1)} \tag{S21} \\
 &= \delta_{i_1 i_2} \cdots \delta_{i_{2m-1} i_{2m}} \sum_{k=0}^m \left[C_b^{(\ell+1)} \right]^{m-k} \left[C_W^{(\ell+1)} \right]^k \\
 &\quad \times \left\{ \frac{1}{n_\ell^k} \sum_{j_1, \dots, j_k=1}^{n_\ell} H_{j_1 j_1 \dots j_k j_k; \alpha_1 \alpha_2 \dots \alpha_{2k-1} \alpha_{2k}}^{(\ell)} + \left[\binom{m}{k} - 1 \text{ others} \right] \right\} \\
 &\quad + [(2m - 1)!! - 1 \text{ other pairings}] .
 \end{aligned}$$

Studiously disentangling cases with different numbers of repetitions in neuron indices (j_1, \dots, j_k) , we notice that at order $O\left(\frac{1}{n}\right)$, terms without repetition or with only one repetition contribute, finding

$$\begin{aligned}
 & \frac{1}{n_\ell^k} \sum_{j_1, \dots, j_k=1}^{n_\ell} H_{j_1 j_1 \dots j_k j_k; \alpha_1 \alpha_2 \dots \alpha_{2k-1} \alpha_{2k}}^{(\ell)} \tag{S22} \\
 &= \left[\langle \sigma(\tilde{z}_{\alpha_1}) \sigma(\tilde{z}_{\alpha_2}) \rangle_{\tilde{K}^{(\ell)}} \cdots \langle \sigma(\tilde{z}_{\alpha_{2k-1}}) \sigma(\tilde{z}_{\alpha_{2k}}) \rangle_{\tilde{K}^{(\ell)}} \right] \\
 &\quad + \frac{1}{n_\ell} \left\{ \left[\langle \sigma(\tilde{z}_{\alpha_1}) \sigma(\tilde{z}_{\alpha_2}) \sigma(\tilde{z}_{\alpha_3}) \sigma(\tilde{z}_{\alpha_4}) \rangle_{\tilde{K}^{(\ell)}} - \langle \sigma(\tilde{z}_{\alpha_1}) \sigma(\tilde{z}_{\alpha_2}) \rangle_{\tilde{K}^{(\ell)}} \langle \sigma(\tilde{z}_{\alpha_3}) \sigma(\tilde{z}_{\alpha_4}) \rangle_{\tilde{K}^{(\ell)}} \right] \right. \\
 &\quad \quad \times \left[\langle \sigma(\tilde{z}_{\alpha_5}) \sigma(\tilde{z}_{\alpha_6}) \rangle_{\tilde{K}^{(\ell)}} \cdots \langle \sigma(\tilde{z}_{\alpha_{2k-1}}) \sigma(\tilde{z}_{\alpha_{2k}}) \rangle_{\tilde{K}^{(\ell)}} \right] \\
 &\quad \quad \left. + \left[\binom{k}{2} - 1 \text{ others} \right] \right\} \\
 &\quad + \frac{1}{n_{\ell-1}} \langle [\sigma(z_{1;\alpha_1}) \sigma(z_{1;\alpha_2}) \cdots \sigma(z_{k;\alpha_{2k-1}}) \sigma(z_{k;\alpha_{2k}})] \{ \mathcal{O}_S[\mathbf{z}] + \mathcal{O}_V[\mathbf{z}] \} \rangle_{K^{(\ell)}} \\
 &\quad + O\left(\frac{1}{n^2}\right)
 \end{aligned}$$

where we used the inductive hierarchical assumption at the ℓ -th layer, i.e., its consequence (HACK) and denoted a single-neuron random vector $\tilde{\mathbf{z}} = \{\tilde{z}_\alpha\}_{\alpha=1, \dots, N_D}$ and the Gaussian integral with the core kernel $\langle \tilde{z}_{\alpha_1} \tilde{z}_{\alpha_2} \rangle_{\tilde{K}^{(\ell)}} = \tilde{K}_{\alpha_1 \alpha_2}^{(\ell)}$. Plugging in expressions (OS,OV) for operators $\mathcal{O}_{S,V}[\mathbf{z}]$,

$$\begin{aligned}
 & \left[\langle \sigma(z_{1,\alpha_1}) \sigma(z_{1,\alpha_2}) \cdots \sigma(z_{k,\alpha_{2k-1}}) \sigma(z_{k,\alpha_{2k}}) \rangle_{K^{(\ell)}} \mathcal{O}_S[\mathbf{z}] \right]_{K^{(\ell)}} \tag{S23} \\
 &= \frac{1}{2} \sum_{\alpha'_1, \alpha'_2} \tilde{S}_{(\ell)}^{\alpha'_1 \alpha'_2} \left\{ \left\langle \sigma(\tilde{z}_{\alpha_1}) \sigma(\tilde{z}_{\alpha_2}) \left(\tilde{z}_{\alpha'_1} \tilde{z}_{\alpha'_2} - \tilde{K}_{\alpha'_1 \alpha'_2}^{(\ell)} \right) \right\rangle_{\tilde{K}^{(\ell)}} \right. \\
 &\quad \quad \left. \times \langle \sigma(\tilde{z}_{\alpha_3}) \sigma(\tilde{z}_{\alpha_4}) \rangle_{\tilde{K}^{(\ell)}} \cdots \langle \sigma(\tilde{z}_{\alpha_{2k-1}}) \sigma(\tilde{z}_{\alpha_{2k}}) \rangle_{\tilde{K}^{(\ell)}} + [(k - 1) \text{ others}] \right\}
 \end{aligned}$$

and

$$\begin{aligned}
 & \left\langle \left[\sigma(\mathbf{z}_{1,\alpha_1}) \sigma(\mathbf{z}_{1,\alpha_2}) \cdots \sigma(\mathbf{z}_{k,\alpha_{2k-1}}) \sigma(\mathbf{z}_{k,\alpha_{2k}}) \right] \mathcal{O}_V[\mathbf{z}] \right\rangle_{K^{(\ell)}} \tag{S24} \\
 = & \frac{1}{8} \sum_{\alpha'_1, \alpha'_2, \alpha'_3, \alpha'_4} \tilde{V}_{(\ell)}^{(\alpha'_1 \alpha'_2)(\alpha'_3 \alpha'_4)} \\
 & \times \left\{ \left\langle \sigma(\tilde{\mathbf{z}}_{\alpha_1}) \sigma(\tilde{\mathbf{z}}_{\alpha_2}) \left(\tilde{\mathbf{z}}_{\alpha'_1} \tilde{\mathbf{z}}_{\alpha'_2} \tilde{\mathbf{z}}_{\alpha'_3} \tilde{\mathbf{z}}_{\alpha'_4} - 2\tilde{\mathbf{z}}_{\alpha'_1} \tilde{\mathbf{z}}_{\alpha'_2} \tilde{K}_{\alpha'_3 \alpha'_4}^{(\ell)} - 4\tilde{\mathbf{z}}_{\alpha'_1} \tilde{\mathbf{z}}_{\alpha'_3} \tilde{K}_{\alpha'_2 \alpha'_4}^{(\ell)} \right. \right. \right. \\
 & \quad \left. \left. \left. + \tilde{K}_{\alpha'_1 \alpha'_2}^{(\ell)} \tilde{K}_{\alpha'_3 \alpha'_4}^{(\ell)} + 2\tilde{K}_{\alpha'_1 \alpha'_3}^{(\ell)} \tilde{K}_{\alpha'_2 \alpha'_4}^{(\ell)} \right) \right\rangle_{\tilde{K}^{(\ell)}} \right. \\
 & \quad \left. \times \langle \sigma(\tilde{\mathbf{z}}_{\alpha_3}) \sigma(\tilde{\mathbf{z}}_{\alpha_4}) \rangle_{\tilde{K}^{(\ell)}} \cdots \langle \sigma(\tilde{\mathbf{z}}_{\alpha_{2k-1}}) \sigma(\tilde{\mathbf{z}}_{\alpha_{2k}}) \rangle_{\tilde{K}^{(\ell)}} + [(k-1) \text{ others}] \right\} \\
 + & \frac{1}{4} \sum_{\alpha'_1, \alpha'_2, \alpha'_3, \alpha'_4} \tilde{V}_{(\ell)}^{(\alpha'_1 \alpha'_2)(\alpha'_3 \alpha'_4)} \\
 & \times \left\{ \left\langle \sigma(\tilde{\mathbf{z}}_{\alpha_1}) \sigma(\tilde{\mathbf{z}}_{\alpha_2}) \left(\tilde{\mathbf{z}}_{\alpha'_1} \tilde{\mathbf{z}}_{\alpha'_2} - \tilde{K}_{\alpha'_1 \alpha'_2}^{(\ell)} \right) \right\rangle_{\tilde{K}^{(\ell)}} \left\langle \sigma(\tilde{\mathbf{z}}_{\alpha_3}) \sigma(\tilde{\mathbf{z}}_{\alpha_4}) \left(\tilde{\mathbf{z}}_{\alpha'_3} \tilde{\mathbf{z}}_{\alpha'_4} - \tilde{K}_{\alpha'_3 \alpha'_4}^{(\ell)} \right) \right\rangle_{\tilde{K}^{(\ell)}} \right. \\
 & \quad \left. \times \langle \sigma(\tilde{\mathbf{z}}_{\alpha_5}) \sigma(\tilde{\mathbf{z}}_{\alpha_6}) \rangle_{\tilde{K}^{(\ell)}} \cdots \langle \sigma(\tilde{\mathbf{z}}_{\alpha_{2k-1}}) \sigma(\tilde{\mathbf{z}}_{\alpha_{2k}}) \rangle_{\tilde{K}^{(\ell)}} + \left[\binom{k}{2} - 1 \text{ others} \right] \right\}
 \end{aligned}$$

As special cases, we obtain expressions advertised in the main text to be contained in this Appendix:

$$\begin{aligned}
 \tilde{A}_{\alpha_1 \alpha_2}^{(\ell)} & \equiv \frac{1}{n_\ell} \sum_{j=1}^{n_\ell} H_{jj; \alpha_1 \alpha_2}^{(\ell)} \tag{S25} \\
 = & \langle \sigma(\tilde{\mathbf{z}}_{\alpha_1}) \sigma(\tilde{\mathbf{z}}_{\alpha_2}) \rangle_{\tilde{K}^{(\ell)}} \\
 & + \frac{1}{n_{\ell-1}} \left[\frac{1}{2} \sum_{\alpha'_1, \alpha'_2} S_{(\ell)}^{\alpha'_1 \alpha'_2} \left\langle \sigma(\tilde{\mathbf{z}}_{\alpha_1}) \sigma(\tilde{\mathbf{z}}_{\alpha_2}) \left(\tilde{\mathbf{z}}_{\alpha'_1} \tilde{\mathbf{z}}_{\alpha'_2} - \tilde{K}_{\alpha'_1 \alpha'_2}^{(\ell)} \right) \right\rangle_{\tilde{K}^{(\ell)}} \right. \\
 & \quad + \frac{1}{8} \sum_{\alpha'_1, \alpha'_2, \alpha'_3, \alpha'_4} \tilde{V}_{(\ell)}^{(\alpha'_1 \alpha'_2)(\alpha'_3 \alpha'_4)} \left\langle \sigma(\tilde{\mathbf{z}}_{\alpha_1}) \sigma(\tilde{\mathbf{z}}_{\alpha_2}) \right. \\
 & \quad \times \left(\tilde{\mathbf{z}}_{\alpha'_1} \tilde{\mathbf{z}}_{\alpha'_2} \tilde{\mathbf{z}}_{\alpha'_3} \tilde{\mathbf{z}}_{\alpha'_4} - 2\tilde{\mathbf{z}}_{\alpha'_1} \tilde{\mathbf{z}}_{\alpha'_2} \tilde{K}_{\alpha'_3 \alpha'_4}^{(\ell)} - 4\tilde{\mathbf{z}}_{\alpha'_1} \tilde{\mathbf{z}}_{\alpha'_3} \tilde{K}_{\alpha'_2 \alpha'_4}^{(\ell)} \right. \\
 & \quad \left. \left. \left. + \tilde{K}_{\alpha'_1 \alpha'_2}^{(\ell)} \tilde{K}_{\alpha'_3 \alpha'_4}^{(\ell)} + 2\tilde{K}_{\alpha'_1 \alpha'_3}^{(\ell)} \tilde{K}_{\alpha'_2 \alpha'_4}^{(\ell)} \right) \right\rangle_{\tilde{K}^{(\ell)}} \right] + O\left(\frac{1}{n^2}\right)
 \end{aligned}$$

and

$$\begin{aligned}
 \tilde{B}_{(\alpha_1\alpha_2)(\alpha_3\alpha_4)}^{(\ell)} &\equiv \frac{1}{n_\ell^2} \sum_{j_1, j_2=1}^{n_\ell} \left[H_{j_1 j_1 j_2 j_2; \alpha_1 \alpha_2 \alpha_3 \alpha_4}^{(\ell)} - H_{j_1 j_1; \alpha_1 \alpha_2}^{(\ell)} H_{j_2 j_2; \alpha_3 \alpha_4}^{(\ell)} \right] \quad (\text{S26}) \\
 &= \frac{1}{n_\ell} \left\{ \langle \sigma(\tilde{\mathbf{z}}_{\alpha_1}) \sigma(\tilde{\mathbf{z}}_{\alpha_2}) \sigma(\tilde{\mathbf{z}}_{\alpha_3}) \sigma(\tilde{\mathbf{z}}_{\alpha_4}) \rangle_{\tilde{K}^{(\ell)}} - \langle \sigma(\tilde{\mathbf{z}}_{\alpha_1}) \sigma(\tilde{\mathbf{z}}_{\alpha_2}) \rangle_{\tilde{K}^{(\ell)}} \langle \sigma(\tilde{\mathbf{z}}_{\alpha_3}) \sigma(\tilde{\mathbf{z}}_{\alpha_4}) \rangle_{\tilde{K}^{(\ell)}} \right. \\
 &\quad \left. + \frac{1}{4} \left(\frac{n_\ell}{n_{\ell-1}} \right) \sum_{\alpha'_1, \alpha'_2, \alpha'_3, \alpha'_4} \tilde{V}_{(\ell)}^{(\alpha'_1 \alpha'_2)(\alpha'_3 \alpha'_4)} \left\langle \sigma(\tilde{\mathbf{z}}_{\alpha_1}) \sigma(\tilde{\mathbf{z}}_{\alpha_2}) (\tilde{\mathbf{z}}_{\alpha'_1} \tilde{\mathbf{z}}_{\alpha'_2} - \tilde{K}_{\alpha'_1 \alpha'_2}^{(\ell)}) \right\rangle_{\tilde{K}^{(\ell)}} \right. \\
 &\quad \left. \times \left\langle \sigma(\tilde{\mathbf{z}}_{\alpha_3}) \sigma(\tilde{\mathbf{z}}_{\alpha_4}) (\tilde{\mathbf{z}}_{\alpha'_3} \tilde{\mathbf{z}}_{\alpha'_4} - \tilde{K}_{\alpha'_3 \alpha'_4}^{(\ell)}) \right\rangle_{\tilde{K}^{(\ell)}} \right\} + O\left(\frac{1}{n^2}\right).
 \end{aligned}$$

Assembling everything,

$$\begin{aligned}
 G_{i_1 \dots i_{2m}; \alpha_1 \dots \alpha_{2m}}^{(\ell+1)} &\quad (\text{S27}) \\
 &= \delta_{i_1 i_2} \cdots \delta_{i_{2m-1} i_{2m}} \prod_{k=1}^m \left[C_b^{(\ell+1)} + C_W^{(\ell+1)} \tilde{A}_{\alpha_{2k-1} \alpha_{2k}}^{(\ell)} \right] \\
 &\quad + [(2m-1)!! - 1 \text{ other pairings}] \\
 &\quad + \delta_{i_1 i_2} \cdots \delta_{i_{2m-1} i_{2m}} \tilde{B}_{(\alpha_1 \alpha_2)(\alpha_3 \alpha_4)}^{(\ell)} \prod_{k=3}^m \left[C_b^{(\ell+1)} + C_W^{(\ell+1)} \tilde{A}_{\alpha_{2k-1} \alpha_{2k}}^{(\ell)} \right] \\
 &\quad + \left\{ \left[3 \times \binom{2m}{4} \times (2m-5)!! - 1 \right] \text{ other } (4, 2, 2, \dots, 2) \text{ clusterings} \right\} \\
 &\quad + O\left(\frac{1}{n^2}\right).
 \end{aligned}$$

In particular,

$$\begin{aligned}
 G_{i_1 i_2; \alpha_1 \alpha_2}^{(\ell+1)} &= \delta_{i_1 i_2} \left[C_b^{(\ell+1)} + C_W^{(\ell+1)} \tilde{A}_{\alpha_1 \alpha_2}^{(\ell)} \right] + O\left(\frac{1}{n^2}\right), \quad (\text{S28}) \\
 G_{i_1 i_2 i_3 i_4; \alpha_1 \alpha_2 \alpha_3 \alpha_4}^{(\ell+1)} \Big|_{\text{connected}} &= \delta_{i_1 i_2} \delta_{i_3 i_4} \tilde{B}_{(\alpha_1 \alpha_2)(\alpha_3 \alpha_4)}^{(\ell)} + \delta_{i_1 i_3} \delta_{i_2 i_4} \tilde{B}_{(\alpha_1 \alpha_3)(\alpha_2 \alpha_4)}^{(\ell)} \\
 &\quad + \delta_{i_1 i_4} \delta_{i_2 i_3} \tilde{B}_{(\alpha_1 \alpha_4)(\alpha_2 \alpha_3)}^{(\ell)} + O\left(\frac{1}{n^2}\right), \quad \text{and} \\
 G_{i_1 i_2 \dots i_{2m-1} i_{2m}; \alpha_1 \alpha_2 \dots \alpha_{2m-1} \alpha_{2m}}^{(\ell+1)} \Big|_{\text{connected}} &= O\left(\frac{1}{n^2}\right), \quad \text{for } 2m \geq 6. \quad (\text{S29})
 \end{aligned}$$

completing our inductive proof. Note that $\tilde{B}_{(\alpha_1 \alpha_2)(\alpha_3 \alpha_4)}^{(\ell)} = O\left(\frac{1}{n}\right)$.

Nowhere in our derivation had we assumed anything about the form of activation functions. The only potential exceptions to our formalism are exponentially growing activation functions – which we never see in practice – that would make the Gaussian integrals unintegrable.

Appendix C. Bestiary of concrete examples

C.1. Quadratic activation

Let us take multilayer perceptrons with quadratic activation, $\sigma(z) = z^2$, and study the distributions of preactivations in the second layer as another illustration of our technology. From the master recursion relations (R1-R3) with the initial condition (R0), Wick's contractions yield

$$\tilde{K}_{\alpha_1\alpha_2}^{(2)} = C_b^{(2)} + C_W^{(2)} \left[\tilde{K}_{\alpha_1\alpha_1}^{(1)} \tilde{K}_{\alpha_2\alpha_2}^{(1)} + 2\tilde{K}_{\alpha_1\alpha_2}^{(1)} \tilde{K}_{\alpha_1\alpha_2}^{(1)} \right], \quad (\text{S30})$$

$$\begin{aligned} \frac{\tilde{V}_{(\alpha_1\alpha_2)(\alpha_3\alpha_4)}^{(2)}}{\left[C_W^{(2)}\right]^2} = & 2 \left[\tilde{K}_{\alpha_1\alpha_1}^{(1)} \tilde{K}_{\alpha_3\alpha_3}^{(1)} \left(\tilde{K}_{\alpha_2\alpha_4}^{(1)}\right)^2 + \tilde{K}_{\alpha_1\alpha_1}^{(1)} \tilde{K}_{\alpha_4\alpha_4}^{(1)} \left(\tilde{K}_{\alpha_2\alpha_3}^{(1)}\right)^2 \right. \\ & \left. + \tilde{K}_{\alpha_2\alpha_2}^{(1)} \tilde{K}_{\alpha_3\alpha_3}^{(1)} \left(\tilde{K}_{\alpha_1\alpha_4}^{(1)}\right)^2 + \tilde{K}_{\alpha_2\alpha_2}^{(1)} \tilde{K}_{\alpha_4\alpha_4}^{(1)} \left(\tilde{K}_{\alpha_1\alpha_3}^{(1)}\right)^2 \right] \\ & + 4 \left[\left(\tilde{K}_{\alpha_1\alpha_3}^{(1)}\right)^2 \left(\tilde{K}_{\alpha_2\alpha_4}^{(1)}\right)^2 + \left(\tilde{K}_{\alpha_1\alpha_4}^{(1)}\right)^2 \left(\tilde{K}_{\alpha_2\alpha_3}^{(1)}\right)^2 \right] \\ & + 8 \left[\tilde{K}_{\alpha_1\alpha_1}^{(1)} \tilde{K}_{\alpha_2\alpha_3}^{(1)} \tilde{K}_{\alpha_3\alpha_4}^{(1)} \tilde{K}_{\alpha_4\alpha_2}^{(1)} + \tilde{K}_{\alpha_2\alpha_2}^{(1)} \tilde{K}_{\alpha_3\alpha_4}^{(1)} \tilde{K}_{\alpha_4\alpha_1}^{(1)} \tilde{K}_{\alpha_1\alpha_3}^{(1)} \right. \\ & \left. + \tilde{K}_{\alpha_3\alpha_3}^{(1)} \tilde{K}_{\alpha_4\alpha_1}^{(1)} \tilde{K}_{\alpha_1\alpha_2}^{(1)} \tilde{K}_{\alpha_2\alpha_4}^{(1)} + \tilde{K}_{\alpha_4\alpha_4}^{(1)} \tilde{K}_{\alpha_1\alpha_2}^{(1)} \tilde{K}_{\alpha_2\alpha_3}^{(1)} \tilde{K}_{\alpha_3\alpha_1}^{(1)} \right] \\ & + 16 \left[\tilde{K}_{\alpha_1\alpha_2}^{(1)} \tilde{K}_{\alpha_1\alpha_3}^{(1)} \tilde{K}_{\alpha_2\alpha_4}^{(1)} \tilde{K}_{\alpha_3\alpha_4}^{(1)} + \tilde{K}_{\alpha_1\alpha_2}^{(1)} \tilde{K}_{\alpha_1\alpha_4}^{(1)} \tilde{K}_{\alpha_2\alpha_3}^{(1)} \tilde{K}_{\alpha_3\alpha_4}^{(1)} \right] \\ & + 16 \tilde{K}_{\alpha_1\alpha_3}^{(1)} \tilde{K}_{\alpha_1\alpha_4}^{(1)} \tilde{K}_{\alpha_2\alpha_3}^{(1)} \tilde{K}_{\alpha_2\alpha_4}^{(1)}, \quad \text{and} \\ \tilde{S}_{\alpha_1\alpha_2}^{(2)} = & 0. \end{aligned} \quad (\text{S32})$$

where $\tilde{K}_{\alpha_1\alpha_2}^{(1)} = C_b^{(1)} + C_W^{(1)} \cdot \left(\frac{\mathbf{x}_{\alpha_1} \cdot \mathbf{x}_{\alpha_2}}{n_0}\right)$. These expressions are used in the main text for the experimental study of finite-width corrections on Bayesian inference.

C.2. Details for single-input cases

The recursive relations simplify drastically for the case of a single input, $N_D = 1$. Setting $C_b^{(\ell)} = 0$ for simplicity and dropping α index, our recursive equations reduce to

$$\tilde{K}^{(\ell+1)} = C_W^{(\ell+1)} \left\langle [\sigma(\tilde{z})]^2 \right\rangle_{\tilde{K}^{(\ell)}}, \quad (\text{S33})$$

$$\frac{\tilde{V}^{(\ell+1)}}{\left(\tilde{K}^{(\ell+1)}\right)^2} = \left(\frac{\left\langle [\sigma(\tilde{z})]^4 \right\rangle_{\tilde{K}^{(\ell)}}}{\left\langle [\sigma(\tilde{z})]^2 \right\rangle_{\tilde{K}^{(\ell)}}^2} - 1 \right) + \frac{1}{4} \left(\frac{n_\ell}{n_{\ell-1}} \right) \left(\frac{\left\langle [\sigma(\tilde{z})]^2 \tilde{z}^2 \right\rangle_{\tilde{K}^{(\ell)}}}{\left\langle [\sigma(\tilde{z})]^2 \right\rangle_{\tilde{K}^{(\ell)}} \tilde{K}^{(\ell)}} - 1 \right)^2 \cdot \frac{\tilde{V}^{(\ell)}}{\left(\tilde{K}^{(\ell)}\right)^2}, \text{ and} \quad (\text{S34})$$

$$\begin{aligned} \frac{\tilde{S}^{(\ell+1)}}{\tilde{K}^{(\ell+1)}} &= \frac{1}{2} \left(\frac{n_\ell}{n_{\ell-1}} \right) \left(\frac{\left\langle [\sigma(\tilde{z})]^2 \tilde{z}^2 \right\rangle_{\tilde{K}^{(\ell)}}}{\left\langle [\sigma(\tilde{z})]^2 \right\rangle_{\tilde{K}^{(\ell)}} \tilde{K}^{(\ell)}} - 1 \right) \cdot \frac{\tilde{S}^{(\ell)}}{\tilde{K}^{(\ell)}} \quad (\text{S35}) \\ &+ \frac{1}{8} \left(\frac{n_\ell}{n_{\ell-1}} \right) \left(\frac{\left\langle [\sigma(\tilde{z})]^2 \tilde{z}^4 \right\rangle_{\tilde{K}^{(\ell)}}}{\left\langle [\sigma(\tilde{z})]^2 \right\rangle_{\tilde{K}^{(\ell)}} \left(\tilde{K}^{(\ell)}\right)^2} - 6 \frac{\left\langle [\sigma(\tilde{z})]^2 \tilde{z}^2 \right\rangle_{\tilde{K}^{(\ell)}}}{\left\langle [\sigma(\tilde{z})]^2 \right\rangle_{\tilde{K}^{(\ell)}} \tilde{K}^{(\ell)}} + 3 \right) \cdot \frac{\tilde{V}^{(\ell)}}{\left(\tilde{K}^{(\ell)}\right)^2}. \end{aligned}$$

C.2.1. MONOMIALS WITH SINGLE INPUT

For monomial activations, $\sigma(z) = z^p$, such as in deep linear networks (Saxe et al., 2013) and quadratic activations (Li et al., 2018),

$$\tilde{K}^{(\ell+1)} = \left[(2p-1)!! C_W^{(\ell+1)} \right] \left(\tilde{K}^{(\ell)} \right)^p, \quad (\text{S36})$$

$$\frac{\tilde{V}^{(\ell+1)}}{\left(\tilde{K}^{(\ell+1)}\right)^2} = \left\{ \frac{(4p-1)!!}{[(2p-1)!!]^2} - 1 \right\} + p^2 \left(\frac{n_\ell}{n_{\ell-1}} \right) \frac{\tilde{V}^{(\ell)}}{\left(\tilde{K}^{(\ell)}\right)^2}, \text{ and} \quad (\text{S37})$$

$$\frac{\tilde{S}^{(\ell+1)}}{\tilde{K}^{(\ell+1)}} = \left(\frac{n_\ell}{n_{\ell-1}} \right) \left[p \frac{\tilde{S}^{(\ell)}}{\tilde{K}^{(\ell)}} + \frac{p(p-1)}{2} \frac{\tilde{V}^{(\ell)}}{\left(\tilde{K}^{(\ell)}\right)^2} \right]. \quad (\text{S38})$$

In particular the four-point vertex solution is given by

$$\frac{1}{n_{\ell-1} p^{2(\ell-1)}} \frac{\tilde{V}^{(\ell)}}{\left(\tilde{K}^{(\ell)}\right)^2} = \left\{ \frac{(4p-1)!!}{[(2p-1)!!]^2} - 1 \right\} \left(\sum_{\ell'=1}^{\ell-1} \frac{1}{n_{\ell'} p^{2\ell'}} \right). \quad (\text{S39})$$

The factor $\left(\sum_{\ell'} \frac{1}{n_{\ell'} p^{2\ell'}} \right)$ generalizes the factor $\left(\sum_{\ell'} \frac{1}{n_{\ell'}} \right)$ for linear and ReLU activations. Following Hanin and Rolnick (2018), this factor guides us to narrow hidden layers as we pass through nonlinear activations for $p > 1$.

C.2.2. RELU WITH SINGLE INPUT

ReLU activation, $\sigma(z) = \max(0, z)$, can also be worked out for a single input through Wick's contractions, noting that the Gaussian integral is halved, yielding

$$\tilde{K}^{(\ell+1)} = \left[\frac{C_W^{(\ell+1)}}{2} \right] \tilde{K}^{(\ell)}, \quad (\text{S40})$$

$$\frac{\tilde{V}^{(\ell+1)}}{\left(\tilde{K}^{(\ell+1)}\right)^2} = 5 + \left(\frac{n_\ell}{n_{\ell-1}}\right) \frac{\tilde{V}^{(\ell)}}{\left(\tilde{K}^{(\ell)}\right)^2}, \quad \text{and} \quad (\text{S41})$$

$$\frac{\tilde{S}^{(\ell+1)}}{\tilde{K}^{(\ell+1)}} = \left(\frac{n_\ell}{n_{\ell-1}}\right) \frac{\tilde{S}^{(\ell)}}{\tilde{K}^{(\ell)}}. \quad (\text{S42})$$

Setting $C_W^{(\ell)} = 2$ for simplicity, these equations can be solved, leading to

$$\tilde{K}^{(\ell)} = \tilde{K}^{(1)} = \frac{\|\mathbf{x}\|_2^2}{n_0}, \quad (\text{S43})$$

$$\frac{1}{n_{\ell-1}} \tilde{V}^{(\ell)} = 5 \left(\sum_{\ell'=1}^{\ell-1} \frac{1}{n_{\ell'}} \right) \left(\tilde{K}^{(1)} \right)^2, \quad \text{and} \quad (\text{S44})$$

$$\tilde{S}^{(\ell)} = 0. \quad (\text{S45})$$

C.3. More experiments on output distributions

Here is an extended version of experiments in Section 4.3. As in the main text, take a single black-white image of hand-written digits from the MNIST dataset as an $n_0 = 784$ -dimensional input, without preprocessing. Set bias variance $C_b^{(\ell)} = 0$, weight variance $C_W^{(\ell)} = C_W$, and use activations $\sigma(z) = z$ (linear) with $C_W = 1$, $\sigma(z) = z^2$ (quadratic) with $C_W = \frac{1}{3}$, and $\sigma(z) = \max(0, z)$ (ReLU) with $C_W = 2$. For all three cases, we consider both depth $L = 2$ with widths $(n_0, n_1, n_2) = (784, n, 1)$ and depth $L = 3$ with widths $(n_0, n_1, n_2, n_3) = (784, n, 2n, 1)$. As in Figure 1, in Figure S1, for each width-parameter n of the hidden layers we record the prior distribution of outputs over 10^6 instances of Gaussian weights and compare it with the theoretical prediction. Results again corroborate our theory.

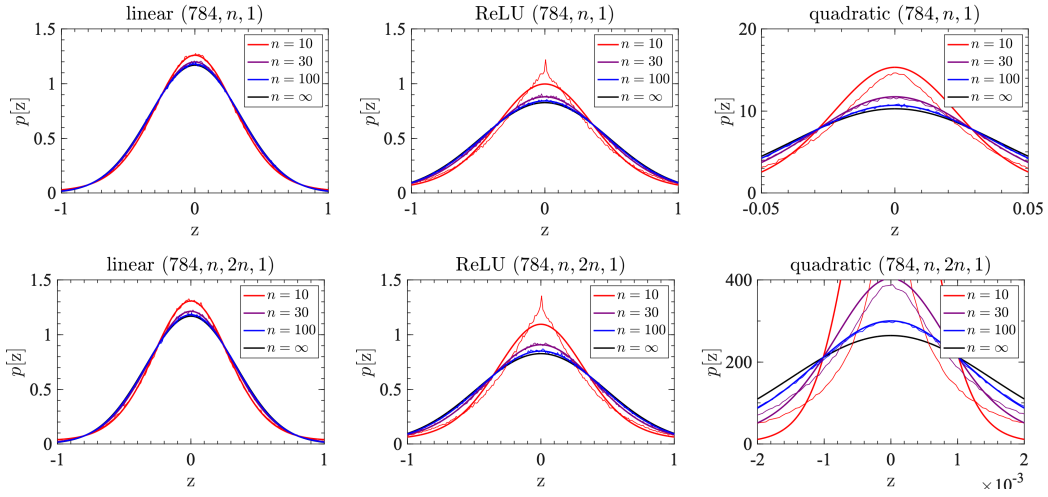


Figure S1: Comparison between theory and experiments for prior distributions of outputs for a single input. Our theoretical predictions (smooth thick lines) and experimental data (rugged thin lines) agree, correctly capturing the initial deviations from the Gaussian processes (black, $n = \infty$), at least down to $n = n_*$ with $n_* \sim 10$ for linear cases, $n_* \sim 30$ for ReLU cases and depth $L = 2$ quadratic case, and $n_* \sim 100$ for depth $L = 3$ quadratic case. This also illustrates that nonlinear activations quickly amplify non-Gaussianity.

Appendix D. Finite-width corrections on Bayesian inference

In order to massage Equation (NGPM) into an actionable form, first playing with the metric inversions and defining $\bar{\phi}_i^\alpha \equiv \sum_{\alpha'} (\tilde{K}^{-1})^{\alpha\alpha'} \bar{\phi}_{i;\alpha'}$, the mean prediction becomes

$$\begin{aligned}
 (y_E^{\text{GP}})_{i;\dot{\gamma}} + \epsilon \sum_{\alpha_1, \dot{\gamma}_1, \alpha_0} (\tilde{K}_\Delta)_{\dot{\gamma}\dot{\gamma}_1} \bar{\phi}_i^{\alpha_1} (\tilde{K}^{-1})^{\dot{\gamma}_1\alpha_0} \\
 \times \left\{ \tilde{S}_{\alpha_0\alpha_1}^{(L)} - \sum_{\alpha_2, \alpha_3} \left[\tilde{V}_{(\alpha_0\alpha_2)(\alpha_1\alpha_3)}^{(L)} + \frac{n_L}{2} \tilde{V}_{(\alpha_0\alpha_1)(\alpha_2\alpha_3)}^{(L)} \right] \right. \\
 \times \left[(\tilde{K}^{-1})^{\alpha_2\alpha_3} - (\tilde{K}^{-1})^{\alpha_2\dot{\gamma}_2} (\tilde{K}_\Delta)_{\dot{\gamma}_2\dot{\gamma}_3} (\tilde{K}^{-1})^{\dot{\gamma}_3\alpha_3} \right] \\
 \left. + \frac{1}{2} \sum_{\alpha_2, \alpha_3} \tilde{V}_{(\alpha_0\alpha_1)(\alpha_2\alpha_3)}^{(L)} \left(\sum_j \bar{\phi}_j^{\alpha_2} \bar{\phi}_j^{\alpha_3} \right) \right\}.
 \end{aligned} \tag{S46}$$

This expression simplifies drastically through the identity

$$\begin{pmatrix} \tilde{K}_{\text{RR}} & \tilde{K}_{\text{RE}} \\ \tilde{K}_{\text{ER}} & \tilde{K}_{\text{EE}} \end{pmatrix}^{-1} = \begin{pmatrix} \tilde{K}_{\text{RR}}^{-1} + \tilde{K}_{\text{RR}}^{-1} \tilde{K}_{\text{RE}} \tilde{K}_\Delta^{-1} \tilde{K}_{\text{ER}} \tilde{K}_{\text{RR}}^{-1} & -\tilde{K}_{\text{RR}}^{-1} \tilde{K}_{\text{RE}} \tilde{K}_\Delta^{-1} \\ -\tilde{K}_\Delta^{-1} \tilde{K}_{\text{ER}} \tilde{K}_{\text{RR}}^{-1} & \tilde{K}_\Delta^{-1} \end{pmatrix}, \tag{S47}$$

which can be checked explicitly, recalling $\tilde{K}_\Delta \equiv \tilde{K}_{\text{EE}} - \tilde{K}_{\text{ER}} \tilde{K}_{\text{RR}}^{-1} \tilde{K}_{\text{RE}}$. Incidentally, this identity can also be used to prove Equation (GPD). Now equipped with this identity, recalling $\bar{\phi}_{i;\alpha} \equiv [(y_{\text{R}})_{i;\bar{\beta}}, (y_{\text{E}}^{\text{GP}})_{i;\dot{\gamma}}]$, we notice that $\bar{\phi}_i^{\bar{\beta}} = \sum_{\bar{\beta}'} (\tilde{K}_{\text{RR}}^{-1})^{\bar{\beta}\bar{\beta}'} (y_{\text{R}})_{i;\bar{\beta}'}$ and $\bar{\phi}_i^{\dot{\gamma}_1} = 0$. Similarly

$$\left[(\tilde{K}^{-1})^{\bar{\beta}_2\bar{\beta}_3} - \sum_{\dot{\gamma}_2, \dot{\gamma}_3} (\tilde{K}^{-1})^{\bar{\beta}_2\dot{\gamma}_2} (\tilde{K}_\Delta)_{\dot{\gamma}_2\dot{\gamma}_3} (\tilde{K}^{-1})^{\dot{\gamma}_3\bar{\beta}_3} \right] = (\tilde{K}_{\text{RR}}^{-1})^{\bar{\beta}_2\bar{\beta}_3}$$

and other components [i.e. with one or both of training components $(\bar{\beta}_2, \bar{\beta}_3)$ replaced by test components $\dot{\gamma}$] vanish. Equation (S46) thus simplifies to

$$\begin{aligned}
 (y_E^{\text{GP}})_{i;\dot{\gamma}} + \epsilon \sum_{\bar{\beta}_1, \dot{\gamma}_1, \alpha_0} (\tilde{K}_\Delta)_{\dot{\gamma}\dot{\gamma}_1} \bar{\phi}_i^{\bar{\beta}_1} (\tilde{K}^{-1})^{\dot{\gamma}_1\alpha_0} \left[\tilde{S}_{\alpha_0\bar{\beta}_1}^{(L)} + \frac{1}{2} \sum_{\bar{\beta}_2, \bar{\beta}_3} \tilde{V}_{(\alpha_0\bar{\beta}_1)(\bar{\beta}_2\bar{\beta}_3)}^{(L)} \left(\sum_j \bar{\phi}_j^{\bar{\beta}_2} \bar{\phi}_j^{\bar{\beta}_3} \right) \right. \\
 \left. - \sum_{\bar{\beta}_2, \bar{\beta}_3} \left(\tilde{V}_{(\alpha_0\bar{\beta}_2)(\bar{\beta}_1\bar{\beta}_3)}^{(L)} + \frac{n_L}{2} \tilde{V}_{(\alpha_0\bar{\beta}_1)(\bar{\beta}_2\bar{\beta}_3)}^{(L)} \right) (\tilde{K}_{\text{RR}}^{-1})^{\bar{\beta}_2\bar{\beta}_3} \right]
 \end{aligned} \tag{S48}$$

Finally, denoting the matrix inside the parenthesis to be

$$\begin{aligned}
 A_{\alpha_0\bar{\beta}_1} \equiv & S_{\alpha_0\bar{\beta}_1}^{(L)} + \frac{1}{2} \sum_{\bar{\beta}_2, \bar{\beta}_3, \bar{\beta}'_2, \bar{\beta}'_3} \tilde{V}_{(\alpha_0\bar{\beta}_1)(\bar{\beta}_2\bar{\beta}_3)}^{(L)} \left(\sum_j \bar{\phi}_j^{\bar{\beta}_2} \bar{\phi}_j^{\bar{\beta}_3} \right) \\
 & - \sum_{\bar{\beta}_2, \bar{\beta}_3} \left(\tilde{V}_{(\alpha_0\bar{\beta}_2)(\bar{\beta}_1\bar{\beta}_3)}^{(L)} + \frac{n_L}{2} \tilde{V}_{(\alpha_0\bar{\beta}_1)(\bar{\beta}_2\bar{\beta}_3)}^{(L)} \right) (\tilde{K}_{\text{RR}}^{-1})^{\bar{\beta}_2\bar{\beta}_3},
 \end{aligned} \tag{NGPM'}$$

and noticing $\sum_{\hat{\gamma}_1} \left(\tilde{K}_\Delta \right)_{\hat{\gamma}_1 \hat{\gamma}_1} \left(\tilde{K}^{-1} \right)^{\hat{\gamma}_1 \bar{\beta}_0} = - \left(\tilde{K}_{\text{ER}} \tilde{K}_{\text{RR}}^{-1} \right)_{\hat{\gamma}}^{\bar{\beta}_0}$ and $\sum_{\hat{\gamma}_1} \left(\tilde{K}_\Delta \right)_{\hat{\gamma}_1 \hat{\gamma}_1} \left(\tilde{K}^{-1} \right)^{\hat{\gamma}_1 \hat{\gamma}_0} = \delta_{\hat{\gamma}}^{\hat{\gamma}_0}$,

$$\left(y_{\text{E}}^{\text{GP}} \right)_{i; \hat{\gamma}} + \epsilon \sum_{\bar{\beta}_1} \bar{\phi}_i^{\bar{\beta}_1} \left[A_{\hat{\gamma} \bar{\beta}_1} - \sum_{\bar{\beta}_0} \left(\tilde{K}_{\text{ER}} \tilde{K}_{\text{RR}}^{-1} \right)_{\hat{\gamma}}^{\bar{\beta}_0} A_{\bar{\beta}_0 \bar{\beta}_1} \right] \quad (\text{NGPM}'')$$

is the mean prediction. Equations (NGPM') and (NGPM'') with $\bar{\phi}_i^{\bar{\beta}} = \sum_{\bar{\beta}'} \left(\tilde{K}_{\text{RR}}^{-1} \right)^{\bar{\beta} \bar{\beta}'} (y_{\text{R}})_{i; \bar{\beta}'}$ are actionable, i.e., easy to program.