

STATISTICAL CHARACTERIZATION OF DEEP NEURAL NETWORKS AND THEIR SENSITIVITY

Anonymous authors

Paper under double-blind review

ABSTRACT

Despite their ubiquity, it remains an active area of research to fully understand deep neural networks (DNNs) and the reasons of their empirical success. We contribute to this effort by introducing a principled approach to statistically characterize DNNs and their sensitivity. By distinguishing between randomness from input data and from model parameters, we study how central and non-central moments of network activation and sensitivity evolve during propagation. Thereby, we provide novel statistical insights on the hypothesis space of input-output mappings encoded by different architectures. Our approach applies both to fully-connected and convolutional networks and incorporates most ingredients of modern DNNs: rectified linear unit (ReLU) activation, batch normalization, skip connections.¹

1 INTRODUCTION

While the empirical success of deep neural networks is not disputed anymore, a full understanding of these models is not yet achieved (Zhang et al., 2016; Dinh et al., 2017; Neyshabur et al., 2017). As no exception to this rule, advances in the design of neural network architectures have more often come from the relentless race of practical applications rather than by principled approaches. Consequently many common practices and rules of thumb are still awaiting for theoretical validation.

An important obstacle in the characterization of neural networks is the complex interplay of different sources of randomness. In that respect, despite winning successes both theoretically (Poole et al., 2016; Schoenholz et al., 2016; Yang & Schoenholz, 2017; Pennington et al., 2018) and empirically (Pennington et al., 2017; Xiao et al., 2018), the mean-field theory of neural networks fails to distinguish between the randomness from input data and model parameters. As a result, input data is only modeled in the rudimentary case of two correlated signals with Gaussian pre-activation distribution. In Balduzzi et al. (2017) input data is similarly modeled using two correlated signals with typical activation patterns. Another path of research considers input data as a 1-dimensional manifold of evolving length and curvature (Poole et al., 2016; Raghu et al., 2017). All cases are limited in their scope or simplifying assumptions.

In this paper, we introduce a novel approach to statistically characterize deep neural networks and their sensitivity. Only mild assumptions are required and the usual simplifications of infinite width, gaussianity or typical activation patterns are not made. Both fully-connected and convolutional networks are encompassed and the commonly used techniques of batch normalization and skip connections are incorporated. The key of our methodology is to consider statistical moments with respect to input data as random variables which depend on model parameters. By studying how different architecture choices influence the evolution with depth of these moments, we provide statistical insights on the corresponding hypothesis spaces of input-output mappings. Our findings span the topics of pseudo-linearity, exploding sensitivity, exponential and power-law evolution with depth.

2 PROPAGATION

We start by formulating the propagation for neural networks with neither batch normalization nor skip connections, that we refer as *vanilla networks*. The formulation will be slightly adapted in section 6 with *batch-normalized feedforward nets*, and in section 7 with *batch-normalized resnets*.

¹Code to reproduce all results will be made available upon publication.

Clean propagation. Suppose that we are given a random tensorial input $\mathbf{x} \in \mathbb{R}^{n \times \dots \times n \times N_0}$ which is spatially d -dimensional with spatial extent of n in each direction and N_0 channels. We further suppose that this input is not trivially zero such that $\mathbb{E}_{\mathbf{x}, \alpha, c}[\mathbf{x}_{\alpha, c}^2] > 0$,² where α denotes the spatial position and c the channel. This input is fed into a d -dimensional convolutional neural network with periodic boundary conditions and constant spatial extent n .³ For each layer, we denote N_l the number of channels or *width*, K_l the convolutional spatial extent, $\omega^l \in \mathbb{R}^{K_l \times \dots \times K_l \times N_{l-1} \times N_l}$ the weight tensors, $\mathbf{b}^l \in \mathbb{R}^{N_l}$ the biases, and $\mathbf{x}^l, \mathbf{y}^l \in \mathbb{R}^{n \times \dots \times n \times N_l}$ the tensors of post-activations and pre-activations. We further denote ϕ the activation function and we adopt the convention $\mathbf{x}^0 \equiv \mathbf{x}$. The propagation at each layer l is given by $\mathbf{x}^l = \phi(\omega^l * \mathbf{x}^{l-1} + \beta^l)$, where $\beta^l \in \mathbb{R}^{n \times \dots \times n \times N_l}$ is the tensor with repeated version of \mathbf{b}^l at each spatial position. From now on, we refer to the propagated tensor \mathbf{x}^l as the *signal*.

Noisy propagation. Next we suppose that the input signal \mathbf{x} is corrupted by a small *white noise* tensor $\epsilon \in \mathbb{R}^{n \times \dots \times n \times N_0}$ with independent and identically distributed components such that $\mathbb{E}_{\epsilon}[\epsilon_i \epsilon_j] = \sigma_{\epsilon}^2 \delta_{ij}$ ($\sigma_{\epsilon} \ll 1$), with δ_{ij} the Kronecker delta. The noisy signal is propagated in the same neural network and we keep track of the noise corruption with the tensor $\epsilon^l \equiv \Phi^l(\mathbf{x} + \epsilon) - \Phi^l(\mathbf{x})$, where Φ^l is the input-output mapping $\mathbf{x}^l = \Phi^l(\mathbf{x})$. Again with the convention $\epsilon^0 \equiv \epsilon$, the simultaneous propagation of the *clean signal* \mathbf{x}^l and the *noise* ϵ^l is given by

$$\mathbf{y}^l = \omega^l * \mathbf{x}^{l-1} + \beta^l, \quad \mathbf{x}^l = \phi(\mathbf{y}^l), \quad (1)$$

$$\epsilon^l = \omega^l * \epsilon^{l-1} \odot \phi'(\mathbf{y}^l), \quad (2)$$

where \odot denotes the element-wise tensor multiplication and Eq. (2) is obtained by taking the derivative in Eq. (1). As shown by Eq. (2), for given \mathbf{x} the mapping from ϵ to ϵ^l is linear. The noise ϵ^l thus stays centered with respect to ϵ during propagation with $\forall \alpha, c: \mathbb{E}_{\epsilon}[\epsilon_{\alpha, c}^l] = 0$.

To get rid of the dependence on σ_{ϵ} , we introduce the random *sensitivity tensor* as the rescaling of the noise with unit initial variance: $\mathbf{s}^0 \equiv \mathbf{s} \equiv \epsilon / \sigma_{\epsilon}$ and $\mathbf{s}^l \equiv \epsilon^l / \sigma_{\epsilon}$. Due to the linearity of Eq. (2), the sensitivity tensor \mathbf{s}^l is the result of the simultaneous propagation of \mathbf{x}^l and \mathbf{s}^l in Eq. (1) and (2). We also have $\mathbb{E}_{\mathbf{s}}[s_i s_j] = \delta_{ij}$ and $\forall \alpha, c: \mathbb{E}_{\mathbf{s}}[s_{\alpha, c}^l] = 0$. The sensitivity tensor encodes *derivative information* while avoiding the burden of increased dimensionality (see computation in Appendix D.1 for an illustration). This will prove very useful.

Further scope. We restrict our analysis to the ReLU *activation function*: $\phi(\mathbf{y}^l) = \max(\mathbf{y}^l, 0)$. This is partly due to lack of space and partly because ReLU networks are the most widely used in practice. Note however that the formulation with convolutional neural networks does not exclude fully-connected networks, obtained simply as a subcase with $n = 1$.

3 INPUT DATA RANDOMNESS – MOMENTS, NORMALIZED SENSITIVITY

3.1 INPUT DATA RANDOMNESS

To understand the importance of the input data distributions $P_{\mathbf{x}}(\mathbf{x}^l)$ and $P_{\mathbf{x}, \mathbf{s}}(\mathbf{s}^l)$, let us adopt a geometrical perspective on Eq. (1) and Eq. (2) in the context of *classification*. The neural network is fed a noisy point cloud with different classes spread throughout space. Its goal is layer after layer to modify this point cloud in $\mathbb{R}^{n \times \dots \times n \times N_l}$ in order to better group points from the same class and separate points from different classes. This continues until the point cloud reaches the final linear separation. The evolution of the point cloud and its derivative across layers, $P_{\mathbf{x}}(\mathbf{x}^l)$ and $P_{\mathbf{x}, \mathbf{s}}(\mathbf{s}^l)$, is of crucial importance since it characterizes the internal neural network machinery.

Note the distinction between $P_{\mathbf{x}}(\mathbf{x}^l)$, $P_{\mathbf{x}, \mathbf{s}}(\mathbf{s}^l)$ on one side, and $P_{\mathbf{x}, \omega, \beta}(\mathbf{x}^l)$, $P_{\mathbf{x}, \mathbf{s}, \omega, \beta}(\mathbf{s}^l)$ on the other side. As an illustration, consider a neural network which has shrunk its input \mathbf{x} to a point mass distribution $P_{\mathbf{x}}(\mathbf{x}^{l-1}) = \delta_{p_{l-1}}$. For given ω^l and β^l , the propagation of Eq. (1) maps this distribution to another point mass $P_{\mathbf{x}}(\mathbf{x}^l) = \delta_{p_l}$. On the other side, $P_{\mathbf{x}, \omega, \beta}(\mathbf{x}^l)$ has density spreading in the whole ambient space $\mathbb{R}^{n \times \dots \times n \times N_l}$ and misses the *distributional pathology*.

²Whenever α and c are considered as random variables they are supposed uniformly sampled among all spatial positions $\{1, \dots, n\}^d$ and all channels $\{1, \dots, N_l\}$.

³The assumptions of periodic boundary conditions and constant spatial extent n are made for simplicity of the analysis. Possible relaxations are discussed in section C.2.

3.2 MOMENTS

In order to keep track of $P_{\mathbf{x}}(\mathbf{x}^l)$, $P_{\mathbf{x},\mathbf{s}}(\mathbf{s}^l)$, our next challenge is the tensorial structure of \mathbf{x}^l and \mathbf{s}^l . In a similar way as batch normalization, we consider feature maps at different spatial positions as interwoven *sub-signals* of the tensorial *meta-signal*. Statistically, we work at the granularity of sub-signals and we treat equally the randomness of the input tensors \mathbf{x} , \mathbf{s} and the spatial position α . This brings the definition of the *feature map vector* and *centered feature map vector* of \mathbf{x}^l as

$$\mathbf{f}_m(\mathbf{x}^l, \alpha) \equiv \mathbf{x}_{\alpha,:}^l, \quad \hat{\mathbf{f}}_m(\mathbf{x}^l, \alpha) \equiv \mathbf{x}_{\alpha,:}^l - \mathbb{E}_{\mathbf{x},\alpha}[\mathbf{x}_{\alpha,:}^l],$$

where α is uniformly sampled in $\{1, \dots, n\}^d$ and the denotation $\mathbf{f}_m(\mathbf{x}^l, \alpha)$ reminds the randomness of both \mathbf{x}^l and α . For any order p , the *non-central moment* and *central moment* of \mathbf{x}^l per-channel and averaged over channels are defined as

$$\begin{aligned} \nu_{p,c}(\mathbf{x}^l) &\equiv \mathbb{E}_{\mathbf{x},\alpha}[\mathbf{f}_m(\mathbf{x}^l, \alpha)_c^p], & \mu_{p,c}(\mathbf{x}^l) &\equiv \mathbb{E}_{\mathbf{x},\alpha}[\hat{\mathbf{f}}_m(\mathbf{x}^l, \alpha)_c^p], \\ \nu_p(\mathbf{x}^l) &\equiv \mathbb{E}_{\mathbf{x},\alpha,c}[\mathbf{f}_m(\mathbf{x}^l, \alpha)_c^p], & \mu_p(\mathbf{x}^l) &\equiv \mathbb{E}_{\mathbf{x},\alpha,c}[\hat{\mathbf{f}}_m(\mathbf{x}^l, \alpha)_c^p]. \end{aligned}$$

The previous definitions are further extended to any random tensor.

3.3 NORMALIZED SENSITIVITY

Normalized sensitivity. Using the moments of \mathbf{x}^l and \mathbf{s}^l we finally define our key metric for the statistical characterization of neural networks that we refer as the *normalized sensitivity* ζ^l :

$$\zeta^l \equiv \left(\frac{\mu_2(\mathbf{s}^l) \mu_2(\mathbf{x}^0)}{\mu_2(\mathbf{x}^l)} \right)^{1/2}, \quad (3)$$

where $\mu_2(\mathbf{s}^l)$ measures *sensitivity* and $\mu_2(\mathbf{x}^l)$, $\mu_2(\mathbf{x}^0)$ measure *signal informativeness*. Again in the classification task where the goal is to set apart different signals, informativeness is measured by $\mu_2(\mathbf{x}^l)$ since a constant shift applied to all signals is uninformative. This is summarized by the property of ζ^l to measure an expected sensitivity when neural network input and output signals are rescaled to have unit variances (proof and visual illustration in Appendix D.2).

Normalized Sensitivity and Signal-to-Noise. Now let us push further the view of noisy propagation with the terminology of *signal-to-noise ratio* SNR^l and *noise factor* F^l :

$$SNR^l \equiv \left(\frac{\sigma_{\text{signal}}^l}{\sigma_{\text{noise}}^l} \right)^2 = \frac{\mu_2(\mathbf{x}^l)}{\mu_2(\epsilon^l)}, \quad F^l \equiv \frac{SNR^0}{SNR^l} = \frac{\mu_2(\mathbf{s}^l) \mu_2(\mathbf{x}^0)}{\mu_2(\mathbf{x}^l)} = (\zeta^l)^2,$$

where we used the definitions of $\mu_2(\mathbf{x}^l)$, $\mu_2(\epsilon^l)$ and $\mu_2(\epsilon^l) = \sigma_\epsilon^2 \mu_2(\mathbf{s}^l)$ and $\mu_2(\epsilon^0) = \sigma_\epsilon^2$. In logarithmic decibel scale, $SNR_{\text{dB}}^l = SNR_{\text{dB}}^0 - 20 \log_{10} \zeta^l$. The normalized sensitivity ζ^l then directly measures how the neural network degrades ($\zeta^l > 1$) or enhances ($\zeta^l < 1$) the input signal-to-noise ratio. The *factor of noise equivalence* means that neural networks with high ζ^l are essentially *noise amplifiers*.

Normalized Sensitivity and Generalization. The relevance of ζ^l is further supported by several studies relating sensitivity and generalization for fully connected networks (Sokolic et al., 2017; Novak et al., 2018; Philipp & Carbonell, 2018). Notably the coefficient defined in Philipp & Carbonell (2018) is equivalent to the normalized sensitivity ζ^l in the fully-connected case (details on equivalence and reasons for our change of terminology in Appendix D.1).

4 MODEL PARAMETERS RANDOMNESS

We now introduce model parameters as the second source of randomness. We only consider *untrained networks* at initialization. Due to the randomness of the initialization point and its independence from input data distribution, this can be seen as characterizing the hypothesis space of input-output mappings encoded by the architecture. We assume that *initialization is standard*: (i) weights and biases are initialized following He et al. (2015), (ii) when pre-activations are

batch-normalized, scale and shift batch normalization parameters are initialized with ones and zeros respectively.

Our methodology from now on is to consider all moment-related quantities such as $\mu_p(\mathbf{x}^l)$, $\mu_p(\mathbf{s}^l)$, $\nu_p(\mathbf{x}^l)$, $\nu_p(\mathbf{s}^l)$, ζ^l as random variables depending on $(\mathbf{W}^1, \mathbf{b}^1, \dots, \mathbf{W}^l, \mathbf{b}^l)$. We introduce the notation $\Theta^l = (\boldsymbol{\omega}^1, \boldsymbol{\beta}^1, \dots, \boldsymbol{\omega}^l, \boldsymbol{\beta}^l)$ for the full set of parameters and the notation $\theta^l = \Theta^l | \Theta^{l-1}$ for the conditional set of parameters, when $(\boldsymbol{\omega}^l, \boldsymbol{\beta}^l)$ are considered as random and $(\boldsymbol{\omega}^1, \boldsymbol{\beta}^1, \boldsymbol{\omega}^{l-1}, \boldsymbol{\beta}^{l-1})$ as given. Furthermore we denote $\delta\mu_2^l(\mathbf{x})$ the geometric increments such that $\log \delta\mu_2(\mathbf{x}^l) = \log \mu_2(\mathbf{x}^l) - \log \mu_2(\mathbf{x}^{l-1})$.

Evolution with Depth. We are now able to write the evolution with depth of $\mu_2(\mathbf{x}^l)$ as

$$\begin{aligned} \log \mu_2(\mathbf{x}^l) - \log \mu_2(\mathbf{x}^0) &= \sum_k \log \mathbb{E}_{\theta^k} [\delta\mu_2(\mathbf{x}^k)] + \sum_k \mathbb{E}_{\theta^k} [\log \delta\mu_2(\mathbf{x}^k)] - \log \mathbb{E}_{\theta^k} [\delta\mu_2(\mathbf{x}^k)] \\ &\quad + \sum_k \log \delta\mu_2(\mathbf{x}^k) - \mathbb{E}_{\theta^k} [\log \delta\mu_2(\mathbf{x}^k)], \end{aligned} \quad (4)$$

where Eq. (4) is obtained using $\log \mu_2(\mathbf{x}^l) - \log \mu_2(\mathbf{x}^0) = \sum_k \log \delta\mu_2(\mathbf{x}^k)$ and by decomposing each term in the sum with telescoping terms. Let us define $\delta_r\mu_2(\mathbf{x}^k)$ such that $\delta\mu_2(\mathbf{x}^k) = \delta_r\mu_2(\mathbf{x}^k) \mathbb{E}_{\theta^k} [\delta\mu_2(\mathbf{x}^k)]$, and denote \bar{m} , \underline{m} , \underline{s} the three different terms in Eq. (4):

$$\bar{m}[\mu_2(\mathbf{x}^k)] \equiv \log \mathbb{E}_{\theta^k} [\delta\mu_2(\mathbf{x}^k)], \quad (5)$$

$$\underline{m}[\mu_2(\mathbf{x}^k)] \equiv \mathbb{E}_{\theta^k} [\log \delta\mu_2(\mathbf{x}^k)] - \log \mathbb{E}_{\theta^k} [\delta\mu_2(\mathbf{x}^k)] = \mathbb{E}_{\theta^k} [\log \delta_r\mu_2(\mathbf{x}^k)], \quad (6)$$

$$\underline{s}[\mu_2(\mathbf{x}^k)] \equiv \log \delta\mu_2(\mathbf{x}^k) - \mathbb{E}_{\theta^k} [\log \delta\mu_2(\mathbf{x}^k)] = \log \delta_r\mu_2(\mathbf{x}^k) - \mathbb{E}_{\theta^k} [\log \delta_r\mu_2(\mathbf{x}^k)], \quad (7)$$

where Eq. (6) is obtained by entering $\log \mathbb{E}_{\theta^k} [\delta\mu_2(\mathbf{x}^k)]$ into the conditional expectation and the logarithm, and Eq. (7) is obtained using $\mathbb{E}_{\theta^k} [\mathbb{E}_{\theta^k} [\log \delta\mu_2(\mathbf{x}^k)]] = \mathbb{E}_{\theta^k} [\log \delta\mu_2(\mathbf{x}^k)]$.

Discussion. First we note that $\bar{m}[\mu_2(\mathbf{x}^k)]$ and $\underline{m}[\mu_2(\mathbf{x}^k)]$ are random variables which depend on Θ^{k-1} while $\underline{s}[\mu_2(\mathbf{x}^k)]$ is a random variable which depends on Θ^k . We also note that $\underline{m}[\mu_2(\mathbf{x}^k)] < 0$ by log-concavity and that $\underline{s}[\mu_2(\mathbf{x}^k)]$ is centered: $\mathbb{E}_{\Theta^k} [\underline{s}[\mu_2(\mathbf{x}^k)]] = 0$.

Under standard initialization, each channel provides an independent contribution to $\mu_2(\mathbf{x}^k) = \mathbb{E}_{\mathbf{x}, \alpha, c} [\hat{\mathbf{f}}_m(\mathbf{x}^k, \boldsymbol{\alpha})_c^2]$. As a consequence, for large N_k the relative increment $\delta_r\mu_2(\mathbf{x}^k)$ has low expected deviation to 1, meaning with high probability that $|\log \delta_r\mu_2(\mathbf{x}^k)| \ll 1$, $|\underline{m}[\mu_2(\mathbf{x}^k)]| \ll 1$, $|\underline{s}[\mu_2(\mathbf{x}^k)]| \ll 1$. In addition, $\underline{s}[\mu_2(\mathbf{x}^k)]$ is centered and *non-correlated* at different k so its sum scales as \sqrt{l} , whereas the sums of $\bar{m}[\mu_2(\mathbf{x}^k)]$ and $\underline{m}[\mu_2(\mathbf{x}^k)]$ scale as l (see Lemma 10 in Appendix E.1). The term $\underline{s}[\mu_2(\mathbf{x}^k)]$ is thus doubly negligible. In summary, the evolution with depth is dominated by $\bar{m}[\mu_2(\mathbf{x}^k)]$ when this term is non-vanishing and by $\underline{m}[\mu_2(\mathbf{x}^k)]$ otherwise. The exact same analysis can be applied to $\nu_2(\mathbf{x}^l)$ and to sensitivity moments $\nu_2(\mathbf{s}^l)$, $\mu_2(\mathbf{s}^l)$. It can also be applied to the ratio $\mu_2(\mathbf{s}^l) / \mu_2(\mathbf{x}^l)$ and thus to ζ^l .

Further notation. From now on, the geometric increment of any quantity is denoted with δ . The definitions of \bar{m} , \underline{m} and \underline{s} in Eq. (5), (6) and (7) are extended to other central and non-central moments of signal and sensitivity as well as ζ^l with $\bar{m}[\zeta^l] = \frac{1}{2}(\bar{m}[\mu_2(\mathbf{s}^l)] - \bar{m}[\mu_2(\mathbf{x}^l)])$, $\underline{m}[\zeta^l] = \frac{1}{2}(\underline{m}[\mu_2(\mathbf{s}^l)] - \underline{m}[\mu_2(\mathbf{x}^l)])$, $\underline{s}[\zeta^l] = \frac{1}{2}(\underline{s}[\mu_2(\mathbf{s}^l)] - \underline{s}[\mu_2(\mathbf{x}^l)])$.

The approximation of dominating terms such as $\bar{m}[\mu_2(\mathbf{x}^k)]$ in Eq. (4) is denoted with \simeq . To make it precise, we write $a \simeq b$ when $a(1 + \delta_a) = b(1 + \delta_b)$ with $|\delta_a| \ll 1$, $|\delta_b| \ll 1$ with high probability. We write $a \lesssim b$ when $a(1 + \delta_a) \leq b(1 + \delta_b)$ with $|\delta_a| \ll 1$, $|\delta_b| \ll 1$ with high probability. From now on we further assume that the *width is large*. We stress that this assumption is milder than mean-field since dominating terms such as $\bar{m}[\mu_2(\mathbf{x}^k)]$ in Eq. (4) remain random variables and since we do not require the signal \mathbf{x}^l to be Gaussian.

5 VANILLA NETWORKS

We are now fully equipped to statistically characterize neural network architectures. We start by analyzing vanilla networks corresponding to the equations of propagation introduced in section 2.

Theorem 1. Moments of vanilla networks. (proof in Appendix E.2) *Denote A_l the event $\{\nu_2(\mathbf{x}^l) > 0\}$ and A'_l the complementary event $\{\nu_2(\mathbf{x}^l) = 0\} = \{\mathbb{P}_{\mathbf{x}, \alpha, c} [\mathbf{x}_{\alpha, c}^l = 0] = 1\}$. Then:*

- (i) $\prod_{k=1}^l (1 - 2^{-N_k}) \leq \mathbb{P}[A_l] \leq \prod_{k=1}^l (1 - 2^{-K_k^d N_{k-1} N_k})$
- (ii) *There exist positive constants $m_{\min}, m_{\max}, v_{\min}, v_{\max} > 0$ and sequences of random variables $(m_l), (m'_l), (s_l), (s'_l)$ such that under A_l , s_l, s'_l are centered and*
- $$\log \nu_2(\mathbf{x}^l) = -lm_l + \sqrt{l}s_l + \log \nu_2(\mathbf{x}^0), \quad m_{\min} \leq m_l \leq m_{\max}, \quad v_{\min} \leq \text{Var}_{\Theta^l|A_l}[s_l] \leq v_{\max},$$
- $$\log \mu_2(\mathbf{s}^l) = -lm'_l + \sqrt{l}s'_l, \quad m_{\min} \leq m'_l \leq m_{\max}, \quad v_{\min} \leq \text{Var}_{\Theta^l|A_l}[s'_l] \leq v_{\max}.$$

Discussion. (i) is related to the collapse of ReLU networks, which is studied in Lu et al. (2018). Any vanilla ReLU network almost surely collapses to 0 but the number of required layer is exponential in the width N_l . In practice, it is not a real problem.

(ii) implies that moments of \mathbf{x}^l and \mathbf{s}^l can be written $\nu_2(\mathbf{x}^l) = \nu_2(\mathbf{x}^0) \exp(-lm_l + \sqrt{l}s_l)$ and $\mu_2(\mathbf{s}^l) = \exp(-lm'_l + \sqrt{l}s'_l)$. This behaviour comes from the particularity of He et al. (2015) to keep stable $\mathbb{E}_{\Theta^l} \mathbb{E}_{\mathbf{x}, \alpha, c}[(\mathbf{x}_{\alpha, c}^l)^2] = \mathbb{E}_{\Theta^l}[\nu_2(\mathbf{x}^l)]$ and $\mathbb{E}_{\Theta^l} \mathbb{E}_{\mathbf{x}, \alpha, c}[(\mathbf{s}_{\alpha, c}^l)^2] = \mathbb{E}_{\Theta^l}[\mu_2(\mathbf{s}^l)]$ during propagation, which results in vanishing $\overline{m}[\nu_2(\mathbf{x}^l)]$ and $\overline{m}[\mu_2(\mathbf{s}^l)]$. The dominating terms in Eq. (4) are then $\underline{m}[\nu_2(\mathbf{x}^k)] < 0$ and $\underline{m}[\mu_2(\mathbf{s}^k)] < 0$ (see Appendix E.3 for details). The small negative drift and the increasing variance of $\log \nu_2(\mathbf{x}^l)$ and $\log \mu_2(\mathbf{s}^l)$ is clear in Fig. 1d and Fig. 1e. It is also clear that the distribution of $\log \nu_2(\mathbf{x}^l)$ and $\log \mu_2(\mathbf{s}^l)$ is nearly Gaussian and therefore the distribution of $\nu_2(\mathbf{x}^l)$ and $\mu_2(\mathbf{s}^l)$ is nearly lognormal. Since a lognormal distribution $\exp(\mu + \sigma X)$ with $X \sim \mathcal{N}(0, 1)$ has expectation equal to $\exp(\mu + \sigma^2/2)$, the increasing variance of $\log \nu_2(\mathbf{x}^l)$ and $\log \mu_2(\mathbf{s}^l)$ must be compensated by small negative drift in order for their exponential $\nu_2(\mathbf{x}^l)$ and $\mu_2(\mathbf{s}^l)$ to have stable expectations $\mathbb{E}_{\Theta^l}[\nu_2(\mathbf{x}^l)]$ and $\mathbb{E}_{\Theta^l}[\mu_2(\mathbf{s}^l)]$. Note that the diffusion happens in log-space since layer composition amounts to multiplicative random effect in real space.

As a consequence of (ii) and Chebyshev's inequality, conditionally on A_l the variables $\nu_2(\mathbf{x}^l)$ and $\mu_2(\mathbf{s}^l)$ still converge in probability to 0 (proof in Appendix E.4). He et al. (2015) thus manages to stabilize expectations with respect to all realizations Θ^l . However in practice we only see a single realization Θ^l and for large l this leads with high probability to vanishing network signal and sensitivity (i.e. activations and gradients). Note that this is a finite-width effect and the terms $\underline{m}[\nu_2(\mathbf{x}^l)], \underline{m}[\mu_2(\mathbf{s}^l)], \underline{s}[\nu_2(\mathbf{x}^l)], \underline{s}[\mu_2(\mathbf{s}^l)]$ also vanish in the limit of infinite width.

Theorem 2. Normalized Sensitivity increments of vanilla networks. (proof in Appendix F.1) Under A_{l-1} , the dominating term in the evolution of the normalized sensitivity is:

$$\delta \zeta^l \simeq \exp\left(\overline{m}_{\text{vanilla}}[\zeta^l]\right) = \left(1 - \mathbb{E}_{c, \theta^l|A_{l-1}}\left[\frac{\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-})}{\mu_2(\mathbf{x}^{l-1})}\right]\right)^{-1/2}, \quad (8)$$

where $\mathbf{y}^{l,+} = \max(\mathbf{y}^l, 0)$ and $\mathbf{y}^{l,-} = \max(-\mathbf{y}^l, 0)$.

Discussion. An immediate consequence of Theorem 2 is that $\delta \zeta^l \gtrsim 1$, meaning that normalized sensitivity always increases with depth for ReLU vanilla networks. To further understand the behaviour of ζ^l we proceed by contradiction. Suppose that there is an event D with probability $\mathbb{P}[D] > 0$ under which $\log \zeta^l$ has a drift larger than the diffusion. Under D the ratio $\mu_2(\mathbf{x}^l) / \nu_2(\mathbf{x}^l)$ then converges in probability to 0 and the variance $\mu_2(\mathbf{x}^l)$ becomes arbitrary smaller than the average magnitude $\nu_2(\mathbf{x}^l) = \mathbb{E}_{\mathbf{x}, \alpha}[\|\mathbf{x}^l\|_2^2] / N_l$ (proof in Appendix F.2). All inputs \mathbf{x} are then mapped to a very localized region and the distribution of \mathbf{x}^l resembles that of a single point. In turn, this implies that with high probability a given pre-activation channel will have all its values concentrated around either a single positive or a single negative value. This means that with high probability either $\mathbf{y}_{:,c}^{l,+} = 0$ for nearly all \mathbf{x}, α , or $\mathbf{y}_{:,c}^{l,-} = 0$ for nearly all \mathbf{x}, α , and thus $\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-}) / \mu_2(\mathbf{x}^{l-1}) \ll 1$ and $\delta \zeta^l \simeq 1$. This contradicts the presence of the drift larger than the diffusion in $\log \zeta^l$.

The previous argument shows that $\delta \zeta^l \simeq 1$ for large l and small diffusion, i.e. large width. This is further supported by the results of our experiments in Fig. 1a and Fig. 1b. A direct consequence of $\delta \zeta^l \simeq 1$ is that the ratio $\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-}) / \mu_2(\mathbf{x}^{l-1})$ in Eq. (8) is constrained to small values, which leaves two possibilities for the *signal distribution* that we illustrate qualitatively:

- (i) If $\max(\nu_{1,c}(\mathbf{y}^{l,+}), \nu_{1,c}(\mathbf{y}^{l,-}))^2 / \mu_2(\mathbf{x}^{l-1}) \rightarrow 0$, then $\nu_{1,c}(|\mathbf{y}|)^2 / \mu_2(\mathbf{x}^{l-1}) \rightarrow 0$ and first-order standardized moments become ill-defined.

- (ii) If $\min(\nu_{1,c}(\mathbf{y}^{l,+}), \nu_{1,c}(\mathbf{y}^{l,-}))^2 / \mu_2(\mathbf{x}^{l-1}) \rightarrow 0$, then $\mathbf{f}_m(\mathbf{y}^l, \boldsymbol{\alpha})$ becomes concentrated on the semi-line generated by its average vector $(\nu_{1,c}(\mathbf{y}^l))_{1 \leq c \leq N_l}$ (proof in Appendix F.3). In this case, the same pattern of activation is seen with probability one with respect to input data and the neural network becomes *linear*.

The situation (ii) is clearly visible in Fig. 1c. Our analysis provides a novel insight in this previously observed phenomenon of coactivation (Balduzzi et al., 2017). Note that the distributional pathology is severe since a 1-dimensional distribution at layer l implies a 1-dimensional input when taking the perspective of layers $l' > l$.

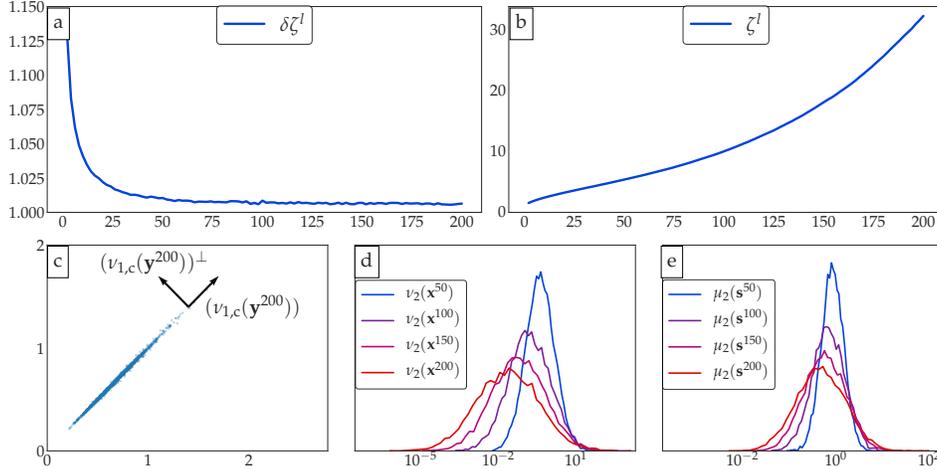


Figure 1: Illustration of the distributional pathologies of vanilla networks with $N_l = 128$ and $L = 200$ layers. (a) Geometric increments of the normalized sensitivity $\delta\zeta^l$ with rapid evolution towards $\delta\zeta^l \simeq 1$. (b) The normalized sensitivity ζ^l has sub-exponential evolution since it is limited by neural network pseudo-linearity. (c) 2-dimensional cuts of preactivation feature maps $\mathbf{f}_m(\mathbf{y}^{200}, \boldsymbol{\alpha})$ using the direction of average vector $(\nu_{1,c}(\mathbf{y}^{200}))$ and a random orthogonal direction $(\nu_{1,c}(\mathbf{y}^{200}))^\perp$. (d-e) Evolution of the distributions of $\mu_2(\mathbf{x}^l)$ and $\mu_2(\mathbf{s}^l)$ with depth showing clear lognormality.

6 BATCH-NORMALIZED FEEDFORWARD NETS

Next we incorporate batch normalization (Ioffe & Szegedy, 2015) which has the effect of subtracting $\nu_{1,c}(\mathbf{y}^l)$ and normalizing by $\mu_{2,c}(\mathbf{y}^l)^{1/2}$ for each channel c in \mathbf{y}^l . The equations of propagation are given by

$$\mathbf{y}^l = \boldsymbol{\omega}^l * \mathbf{x}^{l-1} + \boldsymbol{\beta}^l, \quad \mathbf{z}^l = BN(\mathbf{y}^l), \quad \mathbf{x}^l = \phi(\mathbf{z}^l), \quad (9)$$

$$\mathbf{t}^l = \boldsymbol{\omega}^l * \mathbf{s}^{l-1}, \quad \mathbf{u}^l = BN'(\mathbf{y}^l) \odot \mathbf{t}^l, \quad \mathbf{s}^l = \phi'(\mathbf{z}^l) \odot \mathbf{u}^l, \quad (10)$$

where BN denotes batch normalization and where we introduced the tensors \mathbf{z}^l , \mathbf{t}^l and \mathbf{u}^l . Note that Eq. (9) and Eq. (10) explicitly formulate a finer-grained subdivision of three different steps between layers $l-1$ and l in the simultaneous propagation of $(\mathbf{x}^l, \mathbf{s}^l)$.

Theorem 3. Normalized Sensitivity increments of batch-normalized feedforward nets. (proof in Appendix G.1) *The dominating term in the evolution of ζ^l can be decomposed as the sum of a term $\overline{m}_{BN}[\zeta^l]$ due to batch normalization and a term $\overline{m}_\phi[\zeta^l]$ due to the nonlinearity ϕ :*

$$\exp(\overline{m}_{BN}[\zeta^l]) = \left(\frac{\mu_2(\mathbf{s}^{l-1})}{\mu_2(\mathbf{x}^{l-1})} \right)^{-1/2} \mathbb{E}_{\mathbf{c}, \theta^l} \left[\frac{\mu_{2,c}(\mathbf{t}^l)}{\mu_{2,c}(\mathbf{y}^l)} \right]^{1/2}, \quad (11)$$

$$\exp(\overline{m}_\phi[\zeta^l]) = \left(1 - 2\mathbb{E}_{\mathbf{c}, \theta^l} [\nu_{1,c}(\mathbf{z}^{l,+})\nu_{1,c}(\mathbf{z}^{l,-})] \right)^{-1/2}, \quad (12)$$

$$\delta\zeta^l \simeq \exp(\overline{m}_{BN/FF}[\zeta^l]) = \exp(\overline{m}_{BN}[\zeta^l] + \overline{m}_\phi[\zeta^l]). \quad (13)$$

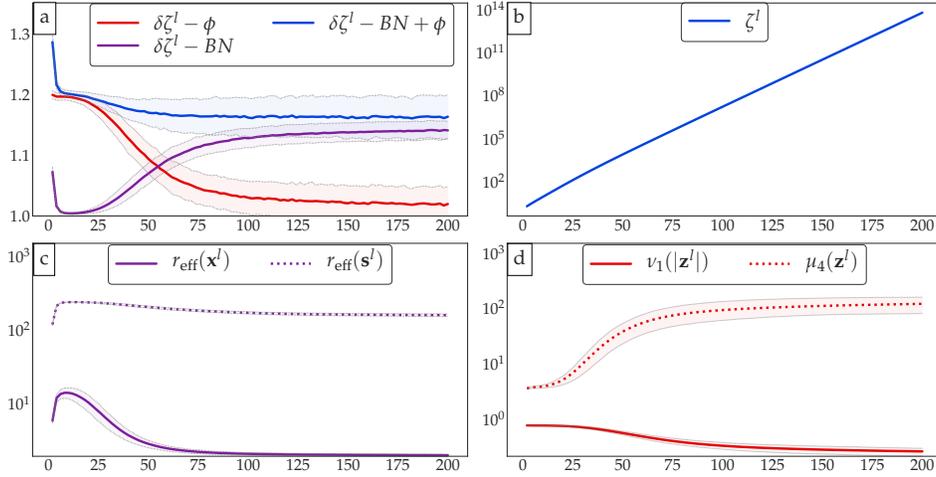


Figure 2: Illustration of the distributional pathology of batch-normalized feedforward nets with $N_l = 384$ and $L = 200$ layers. (a) Geometric increments $\delta\zeta^l$ and their decomposition as the product of two terms: a batch normalization term corresponding to the evolution from $(\mathbf{x}^{l-1}, \mathbf{s}^{l-1})$ to $(\mathbf{z}^l, \mathbf{u}^l)$ and a nonlinearity term corresponding to the evolution from $(\mathbf{z}^l, \mathbf{u}^l)$ to $(\mathbf{x}^l, \mathbf{s}^l)$. (b) The normalized sensitivity ζ^l has exploding behaviour. (c) Effective ranks of signal and sensitivity confirm that sensitivity is much better conditioned than signal. There is a clear inverse correlation between $r_{\text{eff}}(\mathbf{x}^l)$ and the batch normalization term in $\delta\zeta^l$. (d) \mathbf{z}^l becomes ill-behaved as $\nu_1(|\mathbf{z}^l|)$ vanishes and $\mu_4(\mathbf{z}^l)$ explodes. This explains the decay of the nonlinearity term in $\delta\zeta^l$.

Effect of batch normalization. The term of Eq. (11) corresponds to the evolution of ζ^l from $(\mathbf{x}^{l-1}, \mathbf{s}^{l-1})$ at layer $l - 1$ to $(\mathbf{z}^l, \mathbf{u}^l)$ just after *BN*. To understand this term qualitatively, the pre-activation tensor \mathbf{y}^l can be seen as N_l random projections of \mathbf{x}^{l-1} , and batch normalization can be seen as an alteration of the magnitude for each projection. Given that batch normalization uses $\mu_{2,c}(\mathbf{y}^l)^{1/2}$ as normalization factor, directions of high signal variance are dampened while directions of low signal variance are amplified. This *preferential exploration* of low signal directions naturally deteriorates the signal-to-noise ratio and amplifies ζ^l due to the factor of noise equivalence.

Now let us look directly at the quantity inside the expectation in Eq. (11). By spherical symmetry under standard initialization, geometric increments from \mathbf{x}^{l-1} to \mathbf{y}^l for the signal and \mathbf{s}^{l-1} to \mathbf{t}^l for the sensitivity have the same expectation $\mathbb{E}_{c,\theta^l}[\mu_{2,c}(\mathbf{t}^l)] / \mu_{2,c}(\mathbf{s}^{l-1}) = \mathbb{E}_{c,\theta^l}[\mu_{2,c}(\mathbf{y}^l)] / \mu_{2,c}(\mathbf{x}^{l-1})$. On the other hand, the fluctuation of these geometric increments depends on the fluctuation of the signal and sensitivity in the N_l random projections, i.e. on whether directions of signal and sensitivity variances are rare in the ambient space. To measure this effect of *conditioning*, we adopt the metric of effective rank r_{eff} from Vershynin (2012):

$$r_{\text{eff}}(\mathbf{x}^l) \equiv \frac{\text{Tr } \mathbf{C}[\mathbf{f}_m(\mathbf{x}^l, \boldsymbol{\alpha})]}{\|\mathbf{C}[\mathbf{f}_m(\mathbf{x}^l, \boldsymbol{\alpha})]\|}, \quad r_{\text{eff}}(\mathbf{s}^l) \equiv \frac{\text{Tr } \mathbf{C}[\mathbf{f}_m(\mathbf{s}^l, \boldsymbol{\alpha})]}{\|\mathbf{C}[\mathbf{f}_m(\mathbf{s}^l, \boldsymbol{\alpha})]\|}, \quad (14)$$

with $\mathbf{C}[\mathbf{f}_m(\mathbf{x}^l, \boldsymbol{\alpha})]$, $\mathbf{C}[\mathbf{f}_m(\mathbf{s}^l, \boldsymbol{\alpha})]$ the covariance matrices of signal and sensitivity feature map vectors and $\|\cdot\|$ the spectral norm. We further extend this definition to any random tensor. If we assume that \mathbf{s}^l has very good conditioning with $r_{\text{eff}}(\mathbf{s}^{l-1}) \simeq N_{l-1}$, then $\mu_{2,c}(\mathbf{t}^l)$ has small relative deviation to its expectation $\mathbb{E}_{c,\theta^l}[\mu_{2,c}(\mathbf{t}^l)]$ and this term can be treated as a constant. In turn, this implies by convexity of $x \rightarrow 1/x$ that $\exp(\overline{m}_{BN}[\zeta^l]) \gtrsim 1$. The worse the conditioning of \mathbf{x}^{l-1} , i.e. the smaller $r_{\text{eff}}(\mathbf{x}^{l-1})$, the larger the variance of $\mu_{2,c}(\mathbf{y}^l) / \mathbb{E}_{c,\theta^l}[\mu_{2,c}(\mathbf{y}^l)]$ at the denominator and the impact of the convexity. Thus the smaller $r_{\text{eff}}(\mathbf{x}^{l-1})$ and the larger $\exp(\overline{m}_{BN}[\zeta^l])$. This argument is strictly valid for the first step of the propagation where the sensitivity has perfect conditioning, which results in $\exp(\overline{m}_{BN}[\zeta^1]) \geq 1$ (proof in Appendix G.2). In Fig. 2c we confirm experimentally that $r_{\text{eff}}(\mathbf{s}^l) \simeq N_l \gg r_{\text{eff}}(\mathbf{x}^l)$. Together with Fig. 2a we also confirm that $r_{\text{eff}}(\mathbf{x}^l)$ is highly predictive of the batch normalization effect on $\delta\zeta^l$.

Effect of the nonlinearity ϕ . The term of Eq. (12) corresponds to the evolution of ζ^l from $(\mathbf{z}^l, \mathbf{u}^l)$ after BN to $(\mathbf{x}^l, \mathbf{s}^l)$ after ϕ . The quantity inside the expectation can also be expressed as $\nu_{1,c}(\mathbf{z}^{l,+})\nu_{1,c}(\mathbf{z}^{l,-}) / \mu_{2,c}(\mathbf{z}^l)$ since $\mu_{2,c}(\mathbf{z}^l) = 1$ after batch normalization. We then find a very similar expression as Eq. (8) for vanilla networks. The difference is that first-order moments are now normalized using $\mu_{2,c}(\mathbf{z}^l)$ instead of $\mu_2(\mathbf{x}^{l-1})$. This implies that each random projection in \mathbf{y}^l is given the same importance in the sense of similar contribution to $\delta\zeta^l$. On the contrary, Eq. (8) for vanilla networks gives more importance to random directions with high signal since $\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-}) / \mu_2(\mathbf{x}^{l-1})$ is small for low signal directions. Note that $\nu_{1,c}(|\mathbf{z}^l|) = \nu_{1,c}(\mathbf{z}^{l,+}) + \nu_{1,c}(\mathbf{z}^{l,-})$ implies $2\nu_{1,c}(\mathbf{z}^{l,+})\nu_{1,c}(\mathbf{z}^{l,-}) \leq \nu_{1,c}(|\mathbf{z}^l|)^2$ when taking the square. The term $\overline{m}_\phi[\zeta^l]$ is thus limited by $\mathbb{E}_{c,\theta^l}[\nu_{1,c}(|\mathbf{z}^l|)^2]$ as shown by the joint examination of Fig. 2a and Fig. 2d.

7 BATCH-NORMALIZED RESNETS

We finish our exploration of DNN architectures with the incorporation of skip connections. We now suppose that the width is constant, i.e. $N_l = N$ ⁴ and following He et al. (2016) we adopt the perspective of *pre-activation units*. The propagation inside residual units is given by

$$\mathbf{y}^{l,h} = \omega^{l,h} * \phi\left(BN(\mathbf{y}^{l,h-1})\right) + \beta^{l,h}, \quad (15)$$

$$\mathbf{t}^{l,h} = \omega^{l,h} * \left(\mathbf{t}^{l,h-1} \odot BN'(\mathbf{y}^{l,h-1}) \odot \phi'\left(BN(\mathbf{y}^{l,h-1})\right)\right), \quad (16)$$

where Eq. (15) and Eq. (16) hold for $1 \leq h \leq H$, with H the number of layers in each residual unit. Denoting $(\mathbf{y}^{l,h}, \mathbf{t}^{l,h}) = \Phi^{l,h}(\mathbf{y}^{l,h-1}, \mathbf{t}^{l,h-1})$, the propagation between successive residual units is given by

$$(\mathbf{y}^l, \mathbf{t}^l) = (\mathbf{y}^{l-1}, \mathbf{t}^{l-1}) + (\mathbf{y}^{l,H}, \mathbf{t}^{l,H}) = (\mathbf{y}^{l-1}, \mathbf{t}^{l-1}) + \Phi^{l,H} \dots \Phi^{l,1}(\mathbf{y}^{l-1}, \mathbf{t}^{l-1}). \quad (17)$$

For consistency reasons, we rename the inputs of the propagation as $\mathbf{y}^0 \equiv \mathbf{y}$, $\mathbf{t}^0 \equiv \mathbf{t}$. We further adopt the convention that $\mathbf{y}^{0,H} \equiv \mathbf{y}^0$, $\mathbf{t}^{0,H} \equiv \mathbf{t}^0$ such that Eq. (17) can be rewritten as

$$(\mathbf{y}^l, \mathbf{t}^l) = \sum_{k=0}^l (\mathbf{y}^{k,H}, \mathbf{t}^{k,H}). \quad (18)$$

Theorem 4. Normalized Sensitivity increments of batch-normalized resnets. (proof in Appendix H.3) *Suppose that for all depth l we can bound the effective ranks $r_{\min} \lesssim r_{\text{eff}}(\mathbf{y}^l), r_{\text{eff}}(\mathbf{y}^{l,H}), r_{\text{eff}}(\mathbf{t}^l), r_{\text{eff}}(\mathbf{t}^{l,H})$, the second-order central moment $\mu_{2,\min} \lesssim \mu_2(\mathbf{y}^{l,H}) \lesssim \mu_{2,\max}$ and the feedforward increments inside residual units $\delta_{\min} \lesssim \delta\zeta^{l,h} \lesssim \delta_{\max}$. Denote $\rho_{\min} = ((\delta_{\min})^{2H} \mu_{2,\min} - \mu_{2,\max}) / \mu_{2,\max}$ and $\rho_{\max} = ((\delta_{\max})^{2H} \mu_{2,\max} - \mu_{2,\min}) / \mu_{2,\min}$, and further consider τ_{\min}, τ_{\max} such that $\tau_{\min} < \rho_{\min} / 2$ and $\tau_{\max} > \rho_{\max} / 2$.⁵ Then:*

$$\forall l \ll Nr_{\min} : \left(1 + \frac{\rho_{\min}}{l+1}\right)^{1/2} \lesssim \delta\zeta^l \lesssim \left(1 + \frac{\rho_{\max}}{l+1}\right)^{1/2}, \quad (19)$$

$$\forall l \gg 1 : \frac{1}{2}\rho_{\min} \log l \lesssim \log \zeta^l \lesssim \frac{1}{2}\rho_{\max} \log l, \quad (20)$$

$$\forall l \gg 1 : l^{\tau_{\min}} \lesssim \zeta^l \lesssim l^{\tau_{\max}}. \quad (21)$$

Discussion. The evolution in Eq. (19) remains exponential inside residual units since ρ_{\min} and ρ_{\max} have an exponential dependence in H . However it is slowed down by the factor $1/(l+1)$ between successive residual units. This comes from the *dilution* of the residual path $(\mathbf{y}^{l,H}, \mathbf{t}^{l,H})$ in the skip connection path $(\mathbf{y}^{l-1}, \mathbf{t}^{l-1})$ with ratio of signal variances $\mu_2(\mathbf{y}^{l,H}) / \mu_2(\mathbf{y}^{l-1})$ scaling as $1/l$. If we remove the dilution effect in Eq. (19) by replacing $l+1$ by 1 and if we set $\mu_{2,\min} = \mu_{2,\max}$,

⁴Again this assumption is only made for simplicity of the analysis. In practice, it holds at least approximately since N_l is only modified by very few units.

⁵Note that Theorem 4 has very mild assumptions. The assumption on effective ranks is not required in Eq. (20) and Eq. (21). Furthermore the assumption on $\mu_2(\mathbf{y}^{l,H})$ is very reasonable since batch normalization controls signal variance at the beginning of layer H .

then we recover the feedforward evolution with $(\delta_{\min})^H \lesssim \zeta^l \lesssim (\delta_{\max})^H$. The dilution is clearly visible as a side effect of the layer aggregation in Eq. (18): each residual unit l adds a new term of increased sensitivity, but its relative contribution to the aggregation becomes smaller and smaller with l , so it gets harder and harder for the model to grow ζ^l .

Since $\delta\zeta^l$ becomes closer and closer to 1, its fluctuation eventually becomes dominant relatively to its expected deviation to 1. This explains why Eq. (19) only holds for small l . It continues however to hold statistically so that the bounds on $\log \zeta^l = \sum_k \log \delta\zeta^k$ in Eq. (20) correspond to the integration of the bounds in Eq. (19). A direct consequence of the dilution is thus the power-law evolution of ζ^l in Eq. (21) instead of the exponential evolution for feedforward nets. Equivalently, when Eq. (21) is written as $\exp(\tau_{\min} \log l) \lesssim \zeta^l \lesssim \exp(\tau_{\max} \log l)$, the evolution of ζ^l for resnets is the same as the evolution of $\zeta^{\log l}$ for feedforward nets. In words, the evolution with depth of resnets is the *logarithmic* version of the evolution with depth of feedforward nets. Up to some factor, an evolution from 100, to 1 000, and to 10 000 layers for resnets is equivalent to an evolution from 20, to 30, and to 40 layers for feedforward nets. Despite differences in the underlying assumptions, our results are reminiscent of the results in Philipp et al. (2017) on the role of the dilution to alleviate the exploding gradient problem, as well as the results in Yang & Schoenholz (2017) on the power-law evolution of mean-field resnet gradients.

As shown in Fig. 3, the slow power-law evolution ensures that all statistical quantities remain well-behaved. The exponent in the power-law fit in Fig. 3d is set to $\tau = \rho/2 = (\delta_{\text{av}}^{2H} - 1)/2$, with δ_{av} the feedforward increment averaged over the whole evolution. This shows that the bound in Eq. (21) is tight.

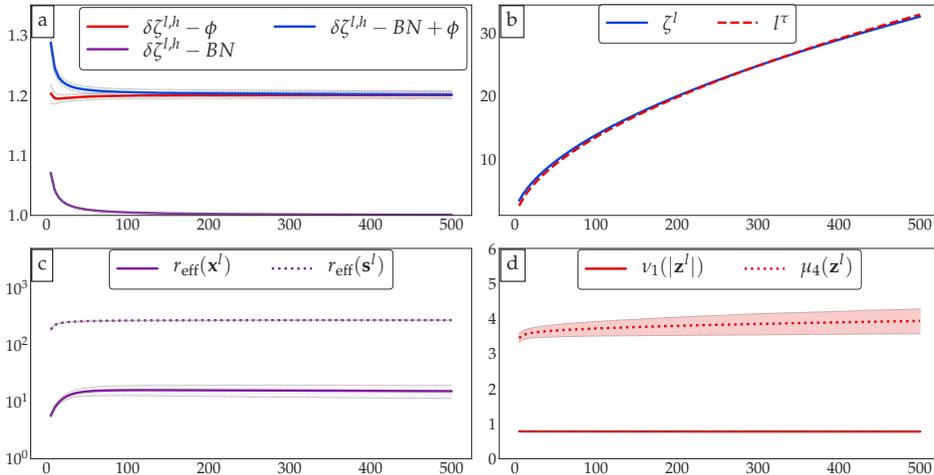


Figure 3: Illustration of the well-behaved evolution of batch-normalized resnets with $N_l = 384$ and $L = 500$ residual units of $H = 2$ layers. (a) Decomposing $\delta\zeta^{l,h}$ as the product of two terms: a batch normalization term and a nonlinearity term. (b) The normalized sensitivity ζ^l has power law evolution. (c) Effective ranks of signal and sensitivity indicate that many directions of signal variance are preserved. (d) Moments $\nu_1(|\mathbf{z}^l|)$ and $\mu_4(\mathbf{z}^l)$ indicate that \mathbf{z}^l has nearly Gaussian distribution.

8 CONCLUSION

This paper introduced a novel approach for the statistical characterization of deep neural networks and their sensitivity. Only very mild assumptions were required and most ingredients of modern DNNs were incorporated. The main scope restriction comes from our focus on the rectifier activation function. We expect however qualitatively similar results to hold for other activation functions.

Below is a summary of our key results:

- For vanilla networks, He et al. (2015) only stabilizes second-order moments of activation and sensitivity in expectation. Depth propagation still induces an additive random walk with small negative drift in log-space. This results in slowly vanishing activations and gradients and the

inevitable convergence to a distributional pathology where the neural network becomes linear and its signal shrunk to a single dimension.

- For batch-normalized feedforward nets, the exponential growth of sensitivity has two origins: on the one hand batch normalization which upweights low-signal pre-activation directions, on the other hand the nonlinear activation function ϕ .
- Finally for resnets the sensitivity only grows as a power-law. Equivalently the evolution with depth of resnets is the logarithmic version of the evolution with depth of feedforward nets. The underlying phenomenon is the dilution of the residual path in the skip connection path with ratio of signal variances decaying as $1/k$. This ingenious mechanism is responsible for breaking the circle of depth multiplicativity which causes distributional pathology for vanilla networks and batch-normalized feedforward nets.

We hope that our methodology will open new perspectives in the statistical understanding of deep neural network architectures. We believe that it could also provide novel insights regarding model trainability and generalization.

REFERENCES

- David Balduzzi, Marcus Frean, Lennox Leary, J. P. Lewis, Kurt Wan-Duo Ma, and Brian McWilliams. The shattered gradients problem: If resnets are the answer, then what is the question? In *International Conference on Machine Learning*, 2017.
- Laurent Dinh, Razvan Pascanu, Samy Bengio, and Yoshua Bengio. Sharp minima can generalize for deep nets. *ArXiv e-prints*, 2017.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the 2015 IEEE International Conference on Computer Vision (ICCV)*, pp. 1026–1034. IEEE Computer Society, 2015.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *Computer Vision – ECCV 2016*, pp. 630–645. Springer International Publishing, 2016.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37, ICML’15*, pp. 448–456, 2015.
- L. Lu, Y. Su, and G. E. Karniadakis. Collapse of Deep and Narrow Neural Nets. *ArXiv e-prints*, 2018.
- Behnam Neyshabur, Srinadh Bhojanapalli, David Mcallester, and Nati Srebro. Exploring generalization in deep learning. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 30*, pp. 5947–5956. 2017.
- R. Novak, Y. Bahri, D. A. Abolafia, J. Pennington, and J. Sohl-Dickstein. Sensitivity and Generalization in Neural Networks: an Empirical Study. In *International Conference on Learning Representations*, 2018.
- Jeffrey Pennington, Samuel S. Schoenholz, and Surya Ganguli. Resurrecting the sigmoid in deep learning through dynamical isometry: theory and practice. *CoRR*, abs/1711.04735, 2017.
- Jeffrey Pennington, Samuel S. Schoenholz, and Surya Ganguli. The emergence of spectral universality in deep networks. In *AISTATS*, volume 84 of *Proceedings of Machine Learning Research*, pp. 1924–1932. PMLR, 2018.
- G. Philipp and J. G. Carbonell. The Nonlinearity Coefficient - Predicting Overfitting in Deep Neural Networks. *ArXiv e-prints*, 2018.
- George Philipp, Dawn Song, and Jaime G. Carbonell. Gradients explode - deep networks are shallow - resnet explained. *ArXiv e-prints*, 2017.

- Ben Poole, Subhaneil Lahiri, Maithra Raghu, Jascha Sohl-Dickstein, and Surya Ganguli. Exponential expressivity in deep neural networks through transient chaos. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 29*, pp. 3360–3368. Curran Associates, Inc., 2016.
- Maithra Raghu, Ben Poole, Jon M. Kleinberg, Surya Ganguli, and Jascha Sohl-Dickstein. On the expressive power of deep neural networks. In *ICML*, volume 70 of *Proceedings of Machine Learning Research*, pp. 2847–2854. PMLR, 2017.
- Samuel S. Schoenholz, Justin Gilmer, Surya Ganguli, and Jascha Sohl-Dickstein. Deep information propagation. *CoRR*, abs/1611.01232, 2016.
- Jure Sokolic, Raja Giryes, Guillermo Sapiro, and Miguel R. D. Rodrigues. Robust large margin deep neural networks. *IEEE Trans. Signal Processing*, 65(16):4265–4280, 2017.
- Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. In *Compressed Sensing: Theory and Applications*, pp. 210–268. Cambridge University Press, 2012.
- L. Xiao, Y. Bahri, J. Sohl-Dickstein, S. S. Schoenholz, and J. Pennington. Dynamical Isometry and a Mean Field Theory of CNNs: How to Train 10,000-Layer Vanilla Convolutional Neural Networks. *ArXiv e-prints*, 2018.
- Ge Yang and Samuel Schoenholz. Mean field residual networks: On the edge of chaos. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 30*, pp. 7103–7114. Curran Associates, Inc., 2017.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. *ArXiv e-prints*, 2016.

A DETAILS OF THE EXPERIMENTS IN FIG. 1, FIG. 2 AND FIG. 3

All three experiments of Fig. 1, Fig. 2 and Fig. 3 were made on CIFAR10 with a random initial convolution of stride 2 reducing the spatial dimension to $n = 16$. In each case we considered the convolutional extent $K_l = 3$ and periodic boundary conditions. A few experiments with approximately statistics-preserving boundary conditions such as *symmetric mirroring* indicated qualitatively equivalent behaviour.

Fig. 1 was obtained by considering width $N_l = 128$ and total depth $L = 200$. For 10 000 random initializations we randomly sampled 1 024 images and computed the evolution with depth of $\nu_2(\mathbf{x}^l)$ and $\mu_2(\mathbf{s}^l)$. The distributions of $\nu_2(\mathbf{x}^l)$ and $\mu_2(\mathbf{s}^l)$ were estimated with the empirical distributions on the 10 000 computations, which are shown in Fig. 1d and Fig. 1e. As for Fig. 1c, it was obtained with 2-dimensional cuts of preactivation feature maps $\mathbf{f}_m(\mathbf{y}^{200}, \boldsymbol{\alpha})$ for 1 randomly sampled point every 5 computations. Each time the cut was performed using the direction of average vector $(\nu_{1,c}(\mathbf{y}^{200}))$ and a random orthogonal direction $(\nu_{1,c}(\mathbf{y}^{200}))^\perp$.

Fig. 2 was obtained by considering 1 000 batches of 64 randomly sampled images. Due to less demanding computations, we increased the width to a more realistic value $N_l = 384$. For simplicity, we always considered batch normalization in train mode. Solid curves correspond to the expectation over the 1 000 batches and Fig. 2a, Fig. 2c, Fig. 2d additionally show 1σ intervals.

Finally Fig. 3 was obtained by considering again 1 000 batches of 64 randomly sampled images with $N_l = 384$ and $L = 500$ residual units of $H = 2$ layers. Geometric increments $\delta\zeta^{l,h}$ were computed in the feedforward propagation from $(\mathbf{y}^{l,h-1}, \mathbf{t}^{l,h-1})$ to $(\mathbf{y}^{l,h}, \mathbf{t}^{l,h})$ at layer $h = 2$ in residual unit l . Solid curves correspond again to the expectation over the 1 000 batches and Fig. 3a, Fig. 3c, Fig. 3d additionally show 1σ intervals.

We plan to release Jupyter notebooks to enable replication of our results upon publication.

B COMPLEMENTARY DEFINITIONS AND NOTATIONS

Receptive field mapping. Here we temporarily need to handle the mechanics of convolution. Let us consider the convolution at layer l of an input $\mathbf{v} \in \mathbb{R}^{n \times \dots \times n \times N_{l-1}}$ from layer $l-1$. The output feature map of the convolution $(w^l * \mathbf{v})_{\alpha,:}$ at position $\alpha \in \{1, \dots, n\}^d$ is obtained by the application of the convolution kernel w^l over a local input region of size $(K_l^d N_{l-1})$, with K_l^d the spatial extent and N_{l-1} the extent in the channel dimension. The local input region is called the *receptive field* of $(w^l * \mathbf{v})$ at spatial position α .

The *receptive field mapping* associates an input \mathbf{v} from layer $l-1$ to $RF(\mathbf{v})$. $RF(\mathbf{v})$ is the tensor of $\mathbb{R}^{n \times \dots \times n \times K_l^d N_{l-1}}$ such that $RF(\mathbf{v})_{\alpha,:}$ is the reshaped vectorial form of the receptive field of $(w^l * \mathbf{v})$ at spatial position α . We denote $R_l = K_l^d N_{l-1}$ the dimensionality of $RF(\mathbf{v})_{\alpha,:}$ and \mathcal{I}_c^l the set of indices in $RF(\mathbf{v})_{\alpha,:}$ corresponding to elements in channel c in the input \mathbf{v} . Strictly speaking, RF depend on l , but this is implied by the argument so we write RF for simplicity.

Receptive field vectors. The *receptive field vector* \mathbf{r}_f and *centered receptive field vector* $\hat{\mathbf{r}}_f$ associated with an input \mathbf{v} from layer l are random vectors which depend on \mathbf{v} , α such that

$$\mathbf{r}_f(\mathbf{v}, \alpha) \equiv RF(\mathbf{v})_{\alpha,:} \quad \text{and} \quad \hat{\mathbf{r}}_f(\mathbf{v}, \alpha) \equiv \mathbf{r}_f(\mathbf{v}, \alpha) - \mathbb{E}_{\mathbf{v}, \alpha}[\mathbf{r}_f(\mathbf{v}, \alpha)],$$

where we kept the same denotation for the variable in the expectation, as a slight abuse of notation. Again \mathbf{r}_f and $\hat{\mathbf{r}}_f$ are strictly speaking dependent on l , but this is implied by the argument.

Statistics-preserving property. RF is *statistics-preserving* with respect to \mathbf{v} if for any channel c and any index $i_c \in \mathcal{I}_c^l$, the random variables $RF(\mathbf{v})_{\alpha, i_c}$ and $\mathbf{v}_{\alpha, c}$ which depend on \mathbf{v} , α have the same distribution $RF(\mathbf{v})_{\alpha, i_c} \stackrel{\mathbf{v}, \alpha}{\sim} \mathbf{v}_{\alpha, c}$.

Equation of Propagation. Using the definition of RF , the affine transformation from the receptive field $RF(\mathbf{x}^{l-1})_{\alpha,:}$ to the feature map in the next layer $\mathbf{y}_{\alpha,:}^l$ is written as

$$\mathbf{y}_{\alpha,:}^l = \mathbf{W}^l RF(\mathbf{x}^{l-1})_{\alpha,:} + \mathbf{b}^l, \quad (22)$$

where $\mathbf{W}^l \in \mathbb{R}^{N_l \times R_l}$ is the suitably reshaped matricial form of ω^l . To lighten notation, we write $\mathbf{y}^l = \mathbf{W}^l RF(\mathbf{x}^{l-1}) + \mathbf{b}^l$ as a short for the affine transformation of Eq. (22) occurring at all spatial positions α . We have the following equivalence between the notations with receptive field and with convolutions:

$$\mathbf{W}^l RF(\mathbf{x}^{l-1}) + \mathbf{b}^l = \omega^l * \mathbf{x}^{l-1} + \beta^l.$$

For vanilla networks, the simultaneous propagation of \mathbf{x}^l and \mathbf{s}^l is then written as

$$\mathbf{y}^l = \mathbf{W}^l RF(\mathbf{x}^{l-1}) + \mathbf{b}^l, \quad \mathbf{x}^l = \phi(\mathbf{y}^l), \quad (23)$$

$$\mathbf{t}^l = \mathbf{W}^l RF(\mathbf{s}^{l-1}), \quad \mathbf{s}^l = \phi'(\mathbf{y}^l) \odot \mathbf{t}^l. \quad (24)$$

For batch-normalized feedforward nets, the simultaneous propagation of \mathbf{x}^l and \mathbf{s}^l is written as

$$\begin{aligned} \mathbf{y}^l &= \mathbf{W}^l RF(\mathbf{x}^{l-1}) + \mathbf{b}^l, & \mathbf{z}^l &= BN(\mathbf{y}^l), & \mathbf{x}^l &= \phi(\mathbf{z}^l), \\ \mathbf{t}^l &= \mathbf{W}^l RF(\mathbf{s}^{l-1}), & \mathbf{u}^l &= BN'(\mathbf{y}^l) \odot \mathbf{t}^l, & \mathbf{s}^l &= \phi'(\mathbf{z}^l) \odot \mathbf{u}^l. \end{aligned}$$

Covariance and Gramian. The *Gramian operator* \mathbf{G} and *covariance operator* \mathbf{C} associated with a random vector $\mathbf{v} \in \mathbb{R}^N$ are defined as

$$\mathbf{G}[\mathbf{v}] \equiv \mathbb{E}_{\mathbf{v}}[\mathbf{v}\mathbf{v}^T] \quad \text{and} \quad \mathbf{C}[\mathbf{v}] \equiv \mathbb{E}_{\mathbf{v}}[\mathbf{v}\mathbf{v}^T] - \mathbb{E}_{\mathbf{v}}[\mathbf{v}]\mathbb{E}_{\mathbf{v}}[\mathbf{v}]^T.$$

Note that the gramian and covariance of feature map and receptive field vectors are related by $\mathbf{G}[\hat{\mathbf{f}}_m(\mathbf{v}, \alpha)] = \mathbf{C}[\hat{\mathbf{f}}_m(\mathbf{v}, \alpha)] = \mathbf{C}[\mathbf{f}_m(\mathbf{v}, \alpha)]$ and $\mathbf{G}[\hat{\mathbf{r}}_f(\mathbf{v}, \alpha)] = \mathbf{C}[\hat{\mathbf{r}}_f(\mathbf{v}, \alpha)] = \mathbf{C}[\mathbf{r}_f(\mathbf{v}, \alpha)]$.

Symmetric propagation for vanilla networks. We define additional tensors obtained by *symmetric propagation* at each layer l . In the case of vanilla networks they are given by:

$$\begin{aligned}\tilde{\mathbf{y}}^l &= -\mathbf{W}^l RF(\mathbf{x}^{l-1}) - \mathbf{b}^l, & \tilde{\mathbf{x}}^l &= \phi(\tilde{\mathbf{y}}^l), \\ \tilde{\mathbf{t}}^l &= -\mathbf{W}^l RF(\mathbf{s}^{l-1}), & \tilde{\mathbf{s}}^l &= \phi'(\tilde{\mathbf{y}}^l) \odot \tilde{\mathbf{t}}^l.\end{aligned}$$

By spherical symmetry, tensor moments have the *same distribution with respect to θ^l for both propagations*. Furthermore $\forall \alpha, c, \mathbf{x}_{\alpha,c}^l + \tilde{\mathbf{x}}_{\alpha,c}^l = |\mathbf{y}_{\alpha,c}^l|$ and $(\mathbf{x}_{\alpha,c}^l)^2 + (\tilde{\mathbf{x}}_{\alpha,c}^l)^2 = (\mathbf{y}_{\alpha,c}^l)^2$ since $\mathbf{x}_{\alpha,c}^l \tilde{\mathbf{x}}_{\alpha,c}^l = 0$. We deduce

$$\forall c : \nu_{2,c}(\mathbf{x}^l) + \nu_{2,c}(\tilde{\mathbf{x}}^l) = \nu_{2,c}(\mathbf{y}^l), \quad (25)$$

Now consider the second-order moments of the sensitivity tensor and suppose first that $\mathbf{x}, \mathbf{s}, \alpha, c, \Theta^{l-1}$ are fixed. We have the following identity:

$$(\mathbf{s}_{\alpha,c}^l)^2 + (\tilde{\mathbf{s}}_{\alpha,c}^l)^2 = (\mathbf{t}_{\alpha,c}^l)^2 \phi'(\mathbf{y}_{\alpha,c}^l)^2 + (\tilde{\mathbf{t}}_{\alpha,c}^l)^2 \phi'(\tilde{\mathbf{y}}_{\alpha,c}^l)^2 = (\mathbf{t}_{\alpha,c}^l)^2 [\phi'(\mathbf{y}_{\alpha,c}^l)^2 + \phi'(\tilde{\mathbf{y}}_{\alpha,c}^l)^2].$$

There are two possible cases depending on $\mathbf{y}_{\alpha,c}^l = \mathbf{W}_{c,:}^l RF(\mathbf{x}^{l-1})_{\alpha,:} + \mathbf{b}^l$:

- If $\|RF(\mathbf{x}^{l-1})_{\alpha,:}\|_2^2 \neq 0$, then under standard initialization, $\mathbb{P}_{\theta^l}[\mathbf{y}_{\alpha,c}^l \neq 0] = 1$, and thus $\mathbb{P}_{\theta^l}[\phi'(\mathbf{y}^l)^2 + \phi'(\tilde{\mathbf{y}}^l)^2 = 1] = 1$ and $\mathbb{P}_{\theta^l}[(\mathbf{s}_{\alpha,c}^l)^2 + (\tilde{\mathbf{s}}_{\alpha,c}^l)^2 = (\mathbf{t}_{\alpha,c}^l)^2] = 1$.
- If $\|RF(\mathbf{x}^{l-1})_{\alpha,:}\|_2^2 = 0$, then the element-wise tensor multiplication of Eq. (24) at layer $l-1$ implies $\|RF(\mathbf{s}^{l-1})_{\alpha,:}\|_2^2 = 0$, and thus $\mathbf{t}_{\alpha,c}^l = 0, \mathbf{s}_{\alpha,c}^l = 0, \tilde{\mathbf{s}}_{\alpha,c}^l = 0$.

In all cases, $\mathbb{P}_{\theta^l}[(\mathbf{s}_{\alpha,c}^l)^2 + (\tilde{\mathbf{s}}_{\alpha,c}^l)^2 = (\mathbf{t}_{\alpha,c}^l)^2] = 1$. For given c , we now consider $\mathbf{x}, \mathbf{s}, \alpha$ as variables again and we get by Fubini's theorem:

$$\begin{aligned}\mathbb{E}_{\theta^l} \mathbb{E}_{\mathbf{x}, \mathbf{s}, \alpha} \left[\mathbf{1}_{\{(\mathbf{s}_{\alpha,c}^l)^2 + (\tilde{\mathbf{s}}_{\alpha,c}^l)^2 = (\mathbf{t}_{\alpha,c}^l)^2\}} \right] &= 1, \\ \mathbb{E}_{\theta^l} \mathbb{E}_{\mathbf{x}, \mathbf{s}, \alpha} \left[\mathbf{1}_{\{(\mathbf{s}_{\alpha,c}^l)^2 + (\tilde{\mathbf{s}}_{\alpha,c}^l)^2 \neq (\mathbf{t}_{\alpha,c}^l)^2\}} \right] &= 0,\end{aligned} \quad (26)$$

$$\mathbb{P}_{\theta^l} \left[\mathbb{P}_{\mathbf{x}, \mathbf{s}, \alpha} [(\mathbf{s}_{\alpha,c}^l)^2 + (\tilde{\mathbf{s}}_{\alpha,c}^l)^2 \neq (\mathbf{t}_{\alpha,c}^l)^2] = 0 \right] = 1. \quad (27)$$

where Eq. (27) is obtained by contradiction, since $\mathbb{P}_{\theta^l} \left[\mathbb{P}_{\mathbf{x}, \mathbf{s}, \alpha} [(\mathbf{s}_{\alpha,c}^l)^2 + (\tilde{\mathbf{s}}_{\alpha,c}^l)^2 \neq (\mathbf{t}_{\alpha,c}^l)^2] > 0 \right] > 0$ would imply that Eq. (26) does not hold. We can relate the moments of $\mathbf{s}^l, \tilde{\mathbf{s}}^l$ and \mathbf{t}^l by

$$\begin{aligned}(\nu_{2,c}(\mathbf{s}^l) + \nu_{2,c}(\tilde{\mathbf{s}}^l) - \nu_{2,c}(\mathbf{t}^l))^2 &= \mathbb{E}_{\mathbf{x}, \mathbf{s}, \alpha} [(\mathbf{s}_{\alpha,c}^l)^2 + (\tilde{\mathbf{s}}_{\alpha,c}^l)^2 - (\mathbf{t}_{\alpha,c}^l)^2]^2 \\ &= \mathbb{E}_{\mathbf{x}, \mathbf{s}, \alpha} \left[\left((\mathbf{s}_{\alpha,c}^l)^2 + (\tilde{\mathbf{s}}_{\alpha,c}^l)^2 - (\mathbf{t}_{\alpha,c}^l)^2 \right) \mathbf{1}_{\{(\mathbf{s}_{\alpha,c}^l)^2 + (\tilde{\mathbf{s}}_{\alpha,c}^l)^2 \neq (\mathbf{t}_{\alpha,c}^l)^2\}} \right]^2 \\ &\leq \nu_{2,c} \left((\mathbf{s}^l)^2 + (\tilde{\mathbf{s}}^l)^2 - (\mathbf{t}^l)^2 \right) \mathbb{P}_{\mathbf{x}, \mathbf{s}, \alpha} [(\mathbf{s}_{\alpha,c}^l)^2 + (\tilde{\mathbf{s}}_{\alpha,c}^l)^2 \neq (\mathbf{t}_{\alpha,c}^l)^2],\end{aligned}$$

where we used Cauchy-Schwarz inequality and the implicit assumption that the moment quantity is well-defined. Combined with Eq. (27), we then get

$$\mathbb{P}_{\theta^l} [\nu_{2,c}(\mathbf{s}^l) + \nu_{2,c}(\tilde{\mathbf{s}}^l) = \nu_{2,c}(\mathbf{t}^l)] = 1.$$

Since $\mathbf{s}^l, \tilde{\mathbf{s}}^l$ and \mathbf{t}^l are centered, $\nu_{2,c}(\mathbf{s}^l) + \nu_{2,c}(\tilde{\mathbf{s}}^l) = \mu_{2,c}(\mathbf{s}^l) + \mu_{2,c}(\tilde{\mathbf{s}}^l)$ and $\nu_{2,c}(\mathbf{t}^l) = \mu_{2,c}(\mathbf{t}^l)$. So we finally get for the event intersection $\{\forall c : \mu_{2,c}(\mathbf{s}^l) + \mu_{2,c}(\tilde{\mathbf{s}}^l) = \mu_{2,c}(\mathbf{t}^l)\} = \bigcap_{c=1}^{N_l} \{\mu_{2,c}(\mathbf{s}^l) + \mu_{2,c}(\tilde{\mathbf{s}}^l) = \mu_{2,c}(\mathbf{t}^l)\}$:

$$\mathbb{P}_{\theta^l} \left[\{\forall c : \mu_{2,c}(\mathbf{s}^l) + \mu_{2,c}(\tilde{\mathbf{s}}^l) = \mu_{2,c}(\mathbf{t}^l)\} \right] = 1. \quad (28)$$

Symmetric propagation for batch-normalized feedforward nets. For batch-normalized feedforward nets, the symmetric propagation at each layer l is given by

$$\tilde{\mathbf{y}}^l = -\mathbf{W}^l RF(\mathbf{x}^{l-1}) - \mathbf{b}^l, \quad \tilde{\mathbf{z}}^l = BN(\tilde{\mathbf{y}}^l), \quad \tilde{\mathbf{x}}^l = \phi(\tilde{\mathbf{z}}^l), \quad (29)$$

$$\tilde{\mathbf{t}}^l = -\mathbf{W}^l RF(\mathbf{s}^{l-1}), \quad \tilde{\mathbf{u}}^l = BN'(\tilde{\mathbf{y}}^l) \odot \tilde{\mathbf{t}}^l, \quad \tilde{\mathbf{s}}^l = \phi'(\tilde{\mathbf{z}}^l) \odot \tilde{\mathbf{u}}^l. \quad (30)$$

BN in Eq. (29) and Eq. (30) use the statistics of $\tilde{\mathbf{y}}^l$, so that tensor moments have the *same distribution with respect to θ^l for both propagations*. We then simply have

$$\tilde{\mathbf{z}}^l = -\mathbf{z}^l, \quad \tilde{\mathbf{u}}^l = -\mathbf{u}^l. \quad (31)$$

The exact same analysis as before gives:

$$\forall \mathbf{c} : \nu_{2,\mathbf{c}}(\mathbf{x}^l) + \nu_{2,\mathbf{c}}(\tilde{\mathbf{x}}^l) = \nu_{2,\mathbf{c}}(\mathbf{z}^l), \quad (32)$$

$$\mathbb{P}_{\theta^l} \left[\left\{ \forall \mathbf{c} : \mu_{2,\mathbf{c}}(\mathbf{s}^l) + \mu_{2,\mathbf{c}}(\tilde{\mathbf{s}}^l) = \mu_{2,\mathbf{c}}(\mathbf{u}^l) \right\} \right] = 1. \quad (33)$$

C STATISTICS-PRESERVING PROPERTY

C.1 CASE OF PERIODIC BOUNDARY CONDITIONS AND CONSTANT SPATIAL EXTENT n

Lemma 1. *If convolutions have periodic boundary conditions and the global spatial extent n is constant, then RF is statistics-preserving with respect to any input \mathbf{v} .*

Proof. Fix a channel \mathbf{c} in \mathbf{v} , an index $i_c \in \mathcal{I}_c^l$, and consider the tensors $\mathbf{v}_{\cdot, \mathbf{c}}$ and $RF(\mathbf{v})_{\cdot, i_c} \in \mathbb{R}^{n \times \dots \times n}$. The index i_c corresponds to a given convolution kernel position $\kappa \in \{1, \dots, K_l\}^d$. Furthermore under periodic boundary conditions, this fixed kernel position implies that each position α in $RF(\mathbf{v})_{\alpha, i_c}$ originates from a different position α' in the tensor $\mathbf{v}_{\alpha', \mathbf{c}}$. Therefore the index mapping $f : \alpha \rightarrow \alpha'$ from $\{1, \dots, n\}^d$ to $\{1, \dots, n\}^d$ is bijective. We then have $RF(\mathbf{v})_{\alpha, i_c} = \mathbf{v}_{f(\alpha), \mathbf{c}} \stackrel{\alpha}{\sim} \mathbf{v}_{\alpha, \mathbf{c}}$ when $RF(\mathbf{v})_{\alpha, i_c}$ and $\mathbf{v}_{\alpha, \mathbf{c}}$ are seen as random variables which depend on α and \mathbf{v} is given. In turn, this implies that $RF(\mathbf{v})_{\alpha, i_c} \stackrel{\mathbf{v}, \alpha}{\sim} \mathbf{v}_{\alpha, \mathbf{c}}$, when they are seen as random variables which depend on \mathbf{v} , α . □

Proposition 2. *If convolutions have periodic boundary conditions and the global spatial extent n is constant, then RF is statistics-preserving with respect to \mathbf{x}^{l-1} and \mathbf{s}^{l-1} .*

Proof. This follows immediately from Lemma 1. □

Corollary 3. *For any \mathbf{c} and $i_c \in \mathcal{I}_c^l$, we have $\mathbf{r}_f(\mathbf{x}^{l-1}, \alpha)_{i_c} \stackrel{\mathbf{x}, \alpha}{\sim} \mathbf{f}_m(\mathbf{x}^{l-1}, \alpha)_{\mathbf{c}}$ and $\mathbf{r}_f(\mathbf{s}^{l-1}, \alpha)_{i_c} \stackrel{\mathbf{x}, \mathbf{s}, \alpha}{\sim} \mathbf{f}_m(\mathbf{s}^{l-1}, \alpha)_{\mathbf{c}}$. Since the cardinality $|\mathcal{I}_c^l| = K_l^d$ is the same for all channels \mathbf{c} , it follows that*

$$\begin{aligned} \nu_2(\mathbf{x}^{l-1}) &= \frac{1}{N_{l-1}} \text{Tr } \mathbf{G}[\mathbf{f}_m(\mathbf{x}^{l-1}, \alpha)] = \frac{1}{R_l} \text{Tr } \mathbf{G}[\mathbf{r}_f(\mathbf{x}^{l-1}, \alpha)], \\ \mu_2(\mathbf{x}^{l-1}) &= \frac{1}{N_{l-1}} \text{Tr } \mathbf{C}[\mathbf{f}_m(\mathbf{x}^{l-1}, \alpha)] = \frac{1}{R_l} \text{Tr } \mathbf{C}[\mathbf{r}_f(\mathbf{x}^{l-1}, \alpha)]. \end{aligned}$$

C.2 RELAXING THE ASSUMPTIONS ON BOUNDARY CONDITIONS AND CONSTANT SPATIAL EXTENT

In this section, we detail possible relaxations of the assumptions on boundary conditions and constant spatial extent n . The global spatial extent is denoted n_l when it is not constant.

C.2.1 CASE OF STATIONARY INPUTS

Periodic extension. The *periodic extension* $\tilde{\mathbf{v}}$ of a random tensor $\mathbf{v} \in \mathbb{R}^{n \times \dots \times n \times N}$ is defined as

$$\tilde{\mathbf{v}}_{\alpha_1+k_1n, \dots, \alpha_d+k_dn, \mathbf{c}} = \mathbf{v}_{\alpha_1, \dots, \alpha_d, \mathbf{c}},$$

with $(k_1, \dots, k_d) \in \mathbb{Z}^d$, and where $\alpha_1, \dots, \alpha_d$ are the d components of the spatial position $\alpha = (\alpha_1, \dots, \alpha_d) \in \{1, \dots, n\}^d$.

Stationarity. The distribution of a random vector $\tilde{\mathbf{v}}$ is defined as *stationary* if, for any k and any configuration of spatial positions $(\alpha_1, \dots, \alpha_k)$ and channels $(\mathbf{c}_1, \dots, \mathbf{c}_k)$, the joint distribution of $\tilde{\mathbf{v}}_{\alpha_1+\alpha_1, \mathbf{c}_1}, \dots, \tilde{\mathbf{v}}_{\alpha_k+\alpha_k, \mathbf{c}_k}$ is the same for all α .

Lemma 4. *If convolutions have periodic boundary conditions and the inputs \mathbf{x} and \mathbf{s} have stationary periodic extensions $\tilde{\mathbf{x}}$ and $\tilde{\mathbf{s}}$, then their periodic extension $\tilde{\mathbf{x}}^l$ and $\tilde{\mathbf{s}}^l$ remain stationary during propagation.*

Proof. First note that convolutions with periodic boundary conditions followed by periodic extensions on \mathbf{v} are equivalent to convolutions on $\tilde{\mathbf{v}}$. This is also the case for componentwise operations such as batch normalization, nonlinear activation and their derivatives. So we can restrict our attention to periodic extensions.

Consider a given spatial shift α and define the translation operator T_α , such that $T_\alpha(\tilde{\mathbf{u}}) = \tilde{\mathbf{v}}$ with $\forall \alpha', c : \tilde{\mathbf{v}}_{\alpha', c} = \tilde{\mathbf{u}}_{\alpha + \alpha', c}$. It is easy to see that T_α commutes with convolutions as well as componentwise operations such as batch normalization, nonlinear activation and their derivatives. It follows that T_α commutes with the input-output mapping Φ^l defined as $(\tilde{\mathbf{x}}^l, \tilde{\mathbf{s}}^l) = \Phi^l(\tilde{\mathbf{x}}, \tilde{\mathbf{s}})$, and thus we have

$$T_\alpha(\tilde{\mathbf{x}}^l, \tilde{\mathbf{s}}^l) = \Phi^l(T_\alpha(\tilde{\mathbf{x}}, \tilde{\mathbf{s}})), \quad (34)$$

where we adopted the notation $T_\alpha(\tilde{\mathbf{u}}, \tilde{\mathbf{v}}) = (T_\alpha(\tilde{\mathbf{u}}), T_\alpha(\tilde{\mathbf{v}}))$. Now consider a given k and configuration of spatial positions $(\alpha_1, \dots, \alpha_k)$ and channels (c_1, \dots, c_k) in $\tilde{\mathbf{x}}^l$ and $\tilde{\mathbf{s}}^l$. Due to limited convolutional spatial extent and due to Eq. (34), there exist a function Φ and a configuration of spatial positions $(\alpha'_1, \dots, \alpha'_{k'})$ and channels $(c'_1, \dots, c'_{k'})$ in $\tilde{\mathbf{x}}$ and $\tilde{\mathbf{s}}$, such that we can write

$$\tilde{\mathbf{x}}_{\alpha_1, c_1}^l, \dots, \tilde{\mathbf{x}}_{\alpha_k, c_k}^l = \Phi(\tilde{\mathbf{x}}_{\alpha'_1, c'_1}, \dots, \tilde{\mathbf{x}}_{\alpha'_{k'}, c'_{k'}}), \quad (35)$$

$$\tilde{\mathbf{x}}_{\alpha + \alpha_1, c_1}^l, \dots, \tilde{\mathbf{x}}_{\alpha + \alpha_k, c_k}^l = \Phi(\tilde{\mathbf{x}}_{\alpha + \alpha'_1, c'_1}, \dots, \tilde{\mathbf{x}}_{\alpha + \alpha'_{k'}, c'_{k'}}), \quad (36)$$

$$\tilde{\mathbf{s}}_{\alpha_1, c_1}^l, \dots, \tilde{\mathbf{s}}_{\alpha_k, c_k}^l = \Phi(\tilde{\mathbf{s}}_{\alpha'_1, c'_1}, \dots, \tilde{\mathbf{s}}_{\alpha'_{k'}, c'_{k'}}), \quad (37)$$

$$\tilde{\mathbf{s}}_{\alpha + \alpha_1, c_1}^l, \dots, \tilde{\mathbf{s}}_{\alpha + \alpha_k, c_k}^l = \Phi(\tilde{\mathbf{s}}_{\alpha + \alpha'_1, c'_1}, \dots, \tilde{\mathbf{s}}_{\alpha + \alpha'_{k'}, c'_{k'}}). \quad (38)$$

By stationarity of $\tilde{\mathbf{x}}$ and $\tilde{\mathbf{s}}$, the terms on the right-hand sides of Eq. (35), Eq. (36) have the same distribution, and the terms on the right-hand sides of Eq. (37), Eq. (38) have the same distribution. It follows that the terms on the left-hand sides of both pairs of equations have the same distribution, meaning that $\tilde{\mathbf{x}}^l$ and $\tilde{\mathbf{s}}^l$ are stationary. \square

Lemma 5. *If convolutions have periodic boundary conditions, then RF is statistics-preserving with respect to any input \mathbf{v} which has stationary periodic extension $\tilde{\mathbf{v}}$ for any global spatial extent n_l .*

Proof. If $\tilde{\mathbf{v}}$ has stationary distribution, it means in particular that for any channel c , the distribution of $\mathbf{v}_{\alpha, c}$ is the same for all $\alpha \in \{1, \dots, n_{l-1}\}^d$. Fix a channel c and an index $i_c \in \mathcal{T}_c^l$ corresponding to a given convolution kernel position $\kappa \in \{1, \dots, K_l\}^d$. A given position α in the receptive field tensor $RF(\mathbf{v})_{\alpha, i_c}$ then corresponds to a given position α' in the original tensor, such that $RF(\mathbf{v})_{\alpha, i_c} = \mathbf{v}_{\alpha', c}$. Since the distribution of $\mathbf{v}_{\alpha', c}$ does not depend on α' , it follows that $RF(\mathbf{v})_{\alpha, i_c} \stackrel{\mathbf{v}; \alpha'}{\approx} \mathbf{v}_{\alpha', c}$ for given α and random α' , and thus that $RF(\mathbf{v})_{\alpha, i_c} \stackrel{\mathbf{v}; \alpha}{\approx} \mathbf{v}_{\alpha, c}$ for random α . \square

Proposition 6. *If convolutions have periodic boundary conditions and the input \mathbf{x} has stationary periodic extension $\tilde{\mathbf{x}}$, then RF is statistics-preserving with respect to \mathbf{x}^l and \mathbf{s}^l , for any global spatial extent n_l .*

Proof. This follows from Lemmas 4 and 5, and from the fact that the input sensitivity tensor has stationary periodic extension $\tilde{\mathbf{s}}$ due to its definition as a white noise tensor with independent and identically distributed components. \square

C.2.2 CASE $n_l \gg K_l$

Proposition 7. *If the convolution stride is one in most layers (i.e. $n_{l-1} = n_l$ in most layers) and the global spatial extent is much larger than the convolutional spatial extent $n_l \gg K_l$ in most layers, then RF is approximately statistics-preserving with respect to \mathbf{x}^{l-1} and \mathbf{s}^{l-1} , for any boundary conditions.*

Proof. Fix a layer $l-1$ such that $n_{l-1} = n_l$ and $n_l \gg K_l$. Denote $RF^{(p)}$ the receptive field mapping at layer l associated with periodic boundary conditions. Since $n_{l-1} = n_l \gg K_l$ the receptive fields $RF(\mathbf{x}^{l-1})_{\alpha,:}$, $RF(\mathbf{s}^{l-1})_{\alpha,:}$ and $RF^{(p)}(\mathbf{x}^{l-1})_{\alpha,:}$, $RF^{(p)}(\mathbf{s}^{l-1})_{\alpha,:}$ do not intersect boundary regions for most α , implying that $RF(\mathbf{x}^{l-1})_{\alpha,:} = RF^{(p)}(\mathbf{x}^{l-1})_{\alpha,:}$, $RF(\mathbf{s}^{l-1})_{\alpha,:} = RF^{(p)}(\mathbf{s}^{l-1})_{\alpha,:}$ for most α . If we denote with F the cumulative distribution functions of any random variable, this implies for any index i_c that $F_{\mathbf{x},\alpha}[RF(\mathbf{x}^{l-1})_{\alpha,i_c}] \simeq F_{\mathbf{x},\alpha}[RF^{(p)}(\mathbf{x}^{l-1})_{\alpha,i_c}]$ and $F_{\mathbf{s},\alpha}[RF(\mathbf{s}^{l-1})_{\alpha,i_c}] \simeq F_{\mathbf{s},\alpha}[RF^{(p)}(\mathbf{s}^{l-1})_{\alpha,i_c}]$.

Since $RF^{(p)}$ is statistics-preserving with respect to \mathbf{x}^{l-1} and \mathbf{s}^{l-1} by Lemma 1, it follows that for any channel c and index $i_c \in \mathcal{I}_c^l$, we have $F_{\mathbf{x},\alpha}[RF^{(p)}(\mathbf{x}^{l-1})_{\alpha,i_c}] = F_{\mathbf{x},\alpha}[\mathbf{x}_{\alpha,c}^{l-1}]$ and $F_{\mathbf{s},\alpha}[RF^{(p)}(\mathbf{s}^{l-1})_{\alpha,i_c}] = F_{\mathbf{s},\alpha}[\mathbf{s}_{\alpha,c}^{l-1}]$. We then deduce that $F_{\mathbf{x},\alpha}[RF(\mathbf{x}^{l-1})_{\alpha,i_c}] \simeq F_{\mathbf{x},\alpha}[\mathbf{x}_{\alpha,c}^{l-1}]$ and $F_{\mathbf{s},\alpha}[RF(\mathbf{s}^{l-1})_{\alpha,i_c}] \simeq F_{\mathbf{s},\alpha}[\mathbf{s}_{\alpha,c}^{l-1}]$, meaning that RF is approximately statistics-preserving for \mathbf{x}^{l-1} and \mathbf{s}^{l-1} . \square

D DEFINITION AND PROPERTIES OF THE NORMALIZED SENSITIVITY

D.1 EQUIVALENCE WITH PREVIOUS DEFINITIONS

In the fully-connected case $n = 1$, Philipp & Carbonell (2018) recently introduced the following coefficient:

$$\left(\frac{\mathbb{E}_{\mathbf{x}} [\|\mathbf{J}^l\|_F^2] \mathbb{E}_c [\text{Var}_{\mathbf{x}}[\mathbf{x}_c^0]]}{N_l \mathbb{E}_c [\text{Var}_{\mathbf{x}}[\mathbf{x}_c^l]]} \right)^{1/2}. \quad (39)$$

Let us prove the equivalence between the definitions of Eq. (3) and Eq. (39). In the fully-connected case, the spatial position α can be ignored and tensors and feature map vectors coincide as $\mathbf{x}^l = \mathbf{f}_m(\mathbf{x}^l)$, $\mathbf{s}^l = \mathbf{f}_m(\mathbf{s}^l) = \hat{\mathbf{f}}_m(\mathbf{s}^l)$. When \mathbf{x} is fixed, the input-output Jacobian $\mathbf{J}^l = \frac{\partial \mathbf{x}^l}{\partial \mathbf{x}^0} \in \mathbb{R}^{N_l \times N_0}$ directly summarizes the propagation of the noise $\epsilon^l = \mathbf{J}^l \epsilon$, and thus the sensitivity $\mathbf{s}^l = \mathbf{J}^l \mathbf{s}$. Due to the white noise property $\mathbb{E}_s [\mathbf{s}_i \mathbf{s}_j] = \delta_{ij}$,

$$\mathbb{E}_{\mathbf{x},\mathbf{s}} \left[\sum_c (\mathbf{s}_c^l)^2 \right] = \mathbb{E}_{\mathbf{x}} [\|\mathbf{J}^l\|_F^2], \quad \mathbb{E}_{\mathbf{x},\mathbf{s},c} [(\mathbf{s}_c^l)^2] = \frac{1}{N_l} \mathbb{E}_{\mathbf{x}} [\|\mathbf{J}^l\|_F^2].$$

Here we clearly see the advantage of the sensitivity tensor to encode information on the Jacobian while avoiding increased dimensionality. Going back to our calculation, the definitions

$$\begin{aligned} \mu_2(\mathbf{s}^l) &= \mathbb{E}_{\mathbf{x},\mathbf{s},c} [\hat{\mathbf{f}}_m(\mathbf{s}^l)_c^2] = \mathbb{E}_{\mathbf{x},\mathbf{s},c} [\mathbf{f}_m(\mathbf{s}^l)_c^2] = \mathbb{E}_{\mathbf{x},\mathbf{s},c} [(\mathbf{s}_c^l)^2], \\ \mu_2(\mathbf{x}^l) &= \mathbb{E}_{\mathbf{x},c} [\hat{\mathbf{f}}_m(\mathbf{x}^l)_c^2] = \mathbb{E}_c [\text{Var}_{\mathbf{x}}[\mathbf{x}_c^l]], \end{aligned}$$

finally give the equivalence between the two definitions:

$$\zeta^l = \left(\frac{\mu_2(\mathbf{s}^l) \mu_2(\mathbf{x}^0)}{\mu_2(\mathbf{x}^l)} \right)^{1/2} = \left(\frac{\mathbb{E}_{\mathbf{x}} [\|\mathbf{J}^l\|_F^2] \mathbb{E}_c [\text{Var}_{\mathbf{x}}[\mathbf{x}_c^0]]}{N_l \mathbb{E}_c [\text{Var}_{\mathbf{x}}[\mathbf{x}_c^l]]} \right)^{1/2}.$$

Philipp & Carbonell (2018) chose the terminology of *nonlinearity coefficient* for this metric. While our analysis unveils a strong relationship between ζ^l and the nonlinearity ϕ , it also reveals a strong relationship with batch normalization which is still a linear operation. So we chose instead the terminology of normalized sensitivity.

D.2 PROPERTY OF NORMALIZED SENSITIVITY

Proposition 8. *The sensitivity tensor and the tensor $\nabla_{\mathbf{x}} \mathbf{x}_{\alpha,c}^l$ containing for given α, c the derivatives of $\mathbf{x}_{\alpha,c}^l$ with respect to \mathbf{x} are related by $\mathbb{E}_s [(\mathbf{s}_{\alpha,c}^l)^2] = \|\nabla_{\mathbf{x}} \mathbf{x}_{\alpha,c}^l\|_2^2$.*

Proof of Proposition 8. Due to the definition of ϵ^l as a small corruption to \mathbf{x}^l , $\epsilon_{\alpha,c}^l$ can be written as a function of $\nabla_{\mathbf{x}} \mathbf{x}_{\alpha,c}^l$ and the input noise ϵ ,

$$\epsilon_{\alpha,c}^l = \langle \nabla_{\mathbf{x}} \mathbf{x}_{\alpha,c}^l, \epsilon \rangle,$$

where $\langle \cdot \rangle$ denotes the standard dot product in input space. It follows from the definition of $\mathbf{s}^l = \epsilon^l / \sigma_{\epsilon}$ that $\mathbf{s}_{\alpha,c}^l = \langle \nabla_{\mathbf{x}} \mathbf{x}_{\alpha,c}^l, \mathbf{s} \rangle$. Due to the white noise property $\mathbb{E}_{\mathbf{s}} [s_i s_j] = \delta_{ij}$, we then get

$$\mathbb{E}_{\mathbf{s}} [(\mathbf{s}_{\alpha,c}^l)^2] = \|\nabla_{\mathbf{x}} \mathbf{x}_{\alpha,c}^l\|_2^2.$$

□

Proposition 9. Define \mathbf{v} the rescaling by a constant factor of \mathbf{x} with unit variance $\mathbb{E}_{\mathbf{x},\alpha,c} [\hat{\mathbf{f}}_m(\mathbf{v}, \alpha)_c^2] = 1$, and \mathbf{v}^l the rescaling by a constant factor of \mathbf{x}^l with unit variance $\mathbb{E}_{\mathbf{x},\alpha,c} [\hat{\mathbf{f}}_m(\mathbf{v}^l, \alpha)_c^2] = 1$. When \mathbf{v}^l is considered as a function of \mathbf{v} , ζ^l measures an expected sensitivity of $\hat{\mathbf{f}}_m(\mathbf{v}^l, \alpha)$ as $\zeta^l = \mathbb{E}_{\mathbf{v},\alpha,c} [\|\nabla_{\mathbf{v}} \hat{\mathbf{f}}_m(\mathbf{v}^l, \alpha)_c\|_2^2]^{1/2}$.

Proof of Proposition 9. First we work with the non-normalized \mathbf{x} and \mathbf{x}^l . We can express the second-order central moment of \mathbf{s} as

$$\mu_2(\mathbf{s}^l) = \mathbb{E}_{\mathbf{x},\alpha,c} \mathbb{E}_{\mathbf{s}} [\mathbf{f}_m(\mathbf{s}^l, \alpha)_c^2] = \mathbb{E}_{\mathbf{x},\alpha,c} \mathbb{E}_{\mathbf{s}} [(\mathbf{s}_{\alpha,c}^l)^2] \quad (40)$$

$$= \mathbb{E}_{\mathbf{x},\alpha,c} [\|\nabla_{\mathbf{x}} \mathbf{x}_{\alpha,c}^l\|_2^2] \quad (41)$$

$$= \mathbb{E}_{\mathbf{x},\alpha,c} [\|\nabla_{\mathbf{x}} \mathbf{f}_m(\mathbf{x}^l, \alpha)_c\|_2^2],$$

where Eq. (40) follows from \mathbf{s}^l being centered and Eq. (41) follows from Proposition 8. Now for given c and α , the difference $\mathbf{f}_m(\mathbf{x}^l, \alpha)_c - \hat{\mathbf{f}}_m(\mathbf{x}^l, \alpha)_c$ is constant with respect to \mathbf{x}^l , and thus with respect to \mathbf{x} . The derivatives of $\mathbf{f}_m(\mathbf{x}^l, \alpha)_c$ and $\hat{\mathbf{f}}_m(\mathbf{x}^l, \alpha)_c$ with respect to \mathbf{x} are then equal for all \mathbf{x} , α and c . Together with the definition of $\mu_2(\mathbf{x}^l)$, we deduce the following:

$$\mu_2(\mathbf{s}^l) = \mathbb{E}_{\mathbf{x},\alpha,c} [\|\nabla_{\mathbf{x}} \hat{\mathbf{f}}_m(\mathbf{x}^l, \alpha)_c\|_2^2],$$

$$\mu_2(\mathbf{x}^l) = \mathbb{E}_{\mathbf{x},\alpha,c} [\hat{\mathbf{f}}_m(\mathbf{x}^l, \alpha)_c^2],$$

$$\zeta^l = \left(\frac{\mu_2(\mathbf{s}^l) \mu_2(\mathbf{x}^0)}{\mu_2(\mathbf{x}^l)} \right)^{1/2} = \left(\frac{\mathbb{E}_{\mathbf{x},\alpha,c} [\|\nabla_{\mathbf{x}} \hat{\mathbf{f}}_m(\mathbf{x}^l, \alpha)_c\|_2^2] \mathbb{E}_{\mathbf{x},\alpha,c} [\hat{\mathbf{f}}_m(\mathbf{x}, \alpha)_c^2]}{\mathbb{E}_{\mathbf{x},\alpha,c} [\hat{\mathbf{f}}_m(\mathbf{x}^l, \alpha)_c^2]} \right)^{1/2}. \quad (42)$$

Defining $\mathbf{v}^l = \mathbf{x}^l / \mathbb{E}_{\mathbf{x},\alpha,c} [\hat{\mathbf{f}}_m(\mathbf{x}^l, \alpha)_c^2]^{1/2}$, $\mathbf{v} = \mathbf{x} / \mathbb{E}_{\mathbf{x},\alpha,c} [\hat{\mathbf{f}}_m(\mathbf{x}, \alpha)_c^2]^{1/2}$, we get $\nabla_{\mathbf{x}} \hat{\mathbf{f}}_m(\mathbf{v}^l, \alpha) = \nabla_{\mathbf{x}} \hat{\mathbf{f}}_m(\mathbf{x}^l, \alpha) / \mathbb{E}_{\mathbf{x},\alpha,c} [\hat{\mathbf{f}}_m(\mathbf{x}^l, \alpha)_c^2]^{1/2}$ and $\nabla_{\mathbf{v}} \hat{\mathbf{f}}_m(\mathbf{v}^l, \alpha) = \nabla_{\mathbf{x}} \hat{\mathbf{f}}_m(\mathbf{v}^l, \alpha) \mathbb{E}_{\mathbf{x},\alpha,c} [\hat{\mathbf{f}}_m(\mathbf{x}, \alpha)_c^2]^{1/2}$. It follows that

$$\mathbb{E}_{\mathbf{v},\alpha,c} [\hat{\mathbf{f}}_m(\mathbf{v}^l, \alpha)_c^2] = 1,$$

$$\mathbb{E}_{\mathbf{v},\alpha,c} [\hat{\mathbf{f}}_m(\mathbf{v}, \alpha)_c^2] = 1,$$

$$\mathbb{E}_{\mathbf{v},\alpha,c} [\|\nabla_{\mathbf{v}} \hat{\mathbf{f}}_m(\mathbf{v}^l, \alpha)_c\|_2^2] = \frac{\mathbb{E}_{\mathbf{x},\alpha,c} [\|\nabla_{\mathbf{x}} \hat{\mathbf{f}}_m(\mathbf{x}^l, \alpha)_c\|_2^2] \mathbb{E}_{\mathbf{x},\alpha,c} [\hat{\mathbf{f}}_m(\mathbf{x}, \alpha)_c^2]}{\mathbb{E}_{\mathbf{x},\alpha,c} [\hat{\mathbf{f}}_m(\mathbf{x}^l, \alpha)_c^2]}. \quad (43)$$

Finally combining Eq. (42) and Eq. (43), we get $\zeta^l = \mathbb{E}_{\mathbf{v},\alpha,c} [\|\nabla_{\mathbf{v}} \hat{\mathbf{f}}_m(\mathbf{v}^l, \alpha)_c\|_2^2]^{1/2}$.

□

Illustration. Let us consider the case of fully-connected networks with L layers, and with 1-dimensional input $N_0 = 1$ and 1-dimensional output $N_L = 1$. Denote $\mathbf{x}^L = \Phi(\mathbf{x})$ the original input-output mapping and $\mathbf{v}^L = \tilde{\Phi}(\mathbf{v})$ the rescaled input-output mapping. Proposition 9 then simply becomes $\zeta^L = \mathbb{E}_{\mathbf{v}} [\tilde{\Phi}'(\mathbf{v})^2]^{1/2}$. In Fig. 4 we show the result of the propagation $\mathbf{v}^L = \Phi(\mathbf{v})$ of an input \mathbf{v} having a mixture of Gaussians distribution in three different cases:

- In Fig. 4a the input \mathbf{v} is propagated through $L = 1$ layer with sigmoid activation. After propagation, the expected derivative is low since each mode of the Gaussian mixture appears in a relatively flat region of the sigmoid activation. This is clearly shown by the input-output data points and by the histogram of the outputs.

- In Fig. 4b the input \mathbf{v} is propagated through $L = 1$ layer with linear activation.
- In Fig. 4c the input \mathbf{v} is propagated through $L = 25$ randomly initialized layers, with ReLU activation in the $L - 1$ first layers and linear activation in the final layer. After propagation, the expected derivative is high due to the erratic behavior of the input-output mapping.

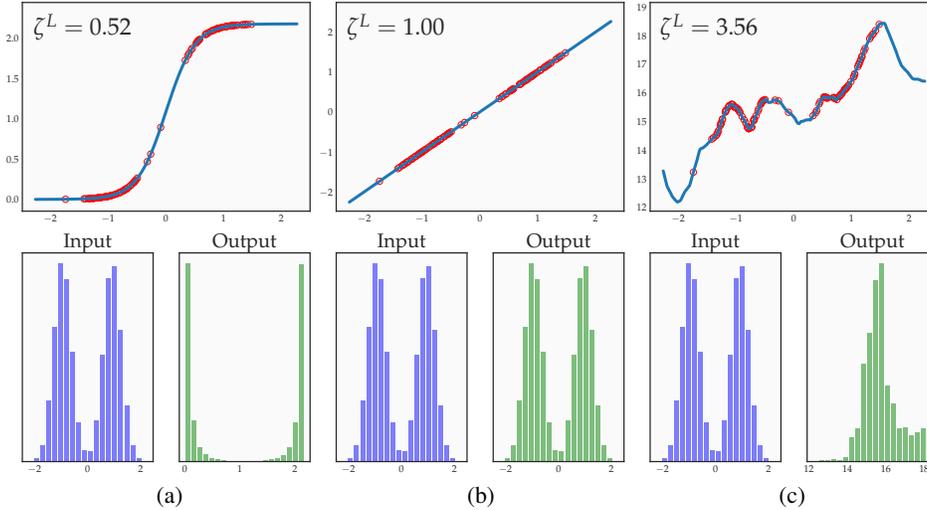


Figure 4: Illustration of the normalized sensitivity for fully-connected networks of L layers, with 1-dimensional input $N_0 = 1$ and 1-dimensional output $N_L = 1$. The distribution of the input \mathbf{v} is a mixture of two Gaussians. We show the result of the propagation in three different cases: (a) $L = 1$ layer with sigmoid activation, (b) $L = 1$ layer with linear activation, (c) $L = 25$ randomly initialized layers, with $N_l = 100$ channels and ReLU activation for $1 \leq l < L$, and linear activation in the final layer $l = L$. *Top*: full input-output mapping (blue curve) and randomly sampled input-output data points (red circles). *Bottom*: histograms of inputs and outputs.

E MOMENTS OF VANILLA NETWORKS

E.1 LEMMA ON THE SUM OF INCREMENTS

Lemma 10. *Let X_k be a sequence of random variables which depend on Θ^k , and denote $Y_k = \mathbb{E}_{\theta^k}[X_k]$ and $Z_k = X_k - \mathbb{E}_{\theta^k}[X_k]$. If there exist constants m_{\min} , m_{\max} , v_{\min} , v_{\max} with $\forall k \leq l$: $m_{\min} \leq Y_k \leq m_{\max}$ and $v_{\min} \leq \text{Var}_{\Theta^k}[Z_k] \leq v_{\max}$, then:*

(i) *The increments Z_k are centered and non-correlated:*

$$\forall k \leq l: \mathbb{E}_{\Theta^k}[Z_k] = 0, \quad \forall k \neq k' \leq l: \mathbb{E}_{\Theta^{\max(k,k')}}[Z_k Z_{k'}] = 0.$$

(ii) *There exist random variables m_l and s_l such that s_l is centered and*

$$\sum_{k=1}^l X_k = l m_l + \sqrt{l} s_l, \quad m_{\min} \leq m_l \leq m_{\max}, \quad v_{\min} \leq \text{Var}_{\Theta^l}[s_l] \leq v_{\max}.$$

Proof of (i). First we show that Z_k is centered:

$$\begin{aligned} \mathbb{E}_{\theta^k}[Z_k] &= \mathbb{E}_{\theta^k}[X_k] - \mathbb{E}_{\theta^k}[X_k] = 0, \\ \mathbb{E}_{\Theta^k}[Z_k] &= \mathbb{E}_{\Theta^{k-1}}[\mathbb{E}_{\theta^k}[Z_k]] = 0. \end{aligned} \tag{44}$$

Now consider $k < k'$, which implies $k \leq k' - 1$ and thus that Z_k is a random variable which only depends on $\Theta^{k'-1}$. Then we can write

$$\begin{aligned}
\mathbb{E}_{\Theta^{k'}}[Z_k Z_{k'}] &= \mathbb{E}_{\Theta^{k'-1}} \mathbb{E}_{\Theta^{k'}}[Z_k Z_{k'}] \\
&= \mathbb{E}_{\Theta^{k'-1}}[Z_k \mathbb{E}_{\Theta^{k'}}[Z_{k'}]] \\
&= 0,
\end{aligned} \tag{45}$$

where Eq. (45) follows from Eq. (44). \square

Proof of (ii). Denote $M_l = \sum_{k=1}^l Y_k$ and $S_l = \sum_{k=1}^l Z_k$. We then have

$$\begin{aligned}
\mathbb{E}_{\Theta^l}[S_l] &= \sum_k \mathbb{E}_{\Theta^l}[Z_k] = 0, \\
\mathbb{E}_{\Theta^l}[S_l^2] &= \sum_{k,k'} \mathbb{E}_{\Theta^l}[Z_k Z_{k'}], \\
\text{Var}_{\Theta^l}[S_l] &= \sum_k \mathbb{E}_{\Theta^k}[Z_k^2] = \sum_k \text{Var}_{\Theta^k}[Z_k],
\end{aligned} \tag{46}$$

where Eq. (46) follows from (i). The hypothesis then gives $lm_{\min} \leq M_l \leq lm_{\max}$ and $lv_{\min} \leq \text{Var}_{\Theta^l}[S_l] \leq lv_{\max}$. If we define $m_l = M_l/l$ and $s_l = S_l/\sqrt{l}$, then s_l is centered and the telescoping sum $\sum_{k=1}^l X_k = \sum_{k=1}^l Y_k + \sum_{k=1}^l Z_k$ can be written as required:

$$\sum_{k=1}^l X_k = M_l + S_l = lm_l + \sqrt{l}s_l, \quad m_{\min} \leq m_l \leq m_{\max}, \quad v_{\min} \leq \text{Var}_{\Theta^l}[s_l] \leq v_{\max}.$$

\square

E.2 PROOF OF THEOREM 1

Theorem 1. Moments of vanilla networks. Denote A_l the event $\{\nu_2(\mathbf{x}^l) > 0\}$ and A'_l the complementary event $\{\nu_2(\mathbf{x}^l) = 0\} = \{\mathbb{P}_{\mathbf{x},\alpha,c}[\mathbf{x}_{\alpha,c}^l = 0] = 1\}$. Then:

- (i) $\prod_{k=1}^l (1 - 2^{-N_k}) \leq \mathbb{P}[A_l] \leq \prod_{k=1}^l (1 - 2^{-K_k^d N_{k-1} N_k})$
- (ii) There exist positive constants $m_{\min}, m_{\max}, v_{\min}, v_{\max} > 0$ and sequences of random variables $(m_l), (m'_l), (s_l), (s'_l)$ such that under A_l, s_l, s'_l are centered and
 - $\log \nu_2(\mathbf{x}^l) = -lm_l + \sqrt{l}s_l + \log \nu_2(\mathbf{x}^0), \quad m_{\min} \leq m_l \leq m_{\max}, \quad v_{\min} \leq \text{Var}_{\Theta^l|A_l}[s_l] \leq v_{\max},$
 - $\log \mu_2(\mathbf{s}^l) = -lm'_l + \sqrt{l}s'_l, \quad m_{\min} \leq m'_l \leq m_{\max}, \quad v_{\min} \leq \text{Var}_{\Theta^l|A_l}[s'_l] \leq v_{\max}.$

Proof of (i). We use the definitions and notations from section B. We further denote $(\mathbf{e}_1, \dots, \mathbf{e}_{R_l})$ and $(\lambda_1, \dots, \lambda_{R_l})$ respectively the orthogonal eigenvectors and eigenvalues of $\mathbf{G}[\mathbf{r}_f(\mathbf{x}^{l-1}, \boldsymbol{\alpha})]$, and $\tilde{\mathbf{W}}^l = \mathbf{W}^l(\mathbf{e}_1, \dots, \mathbf{e}_{R_l})$. We then get

$$\begin{aligned}
\forall c: \nu_{2,c}(\mathbf{y}^l) &= \mathbb{E}_{\mathbf{x},\alpha}[(\mathbf{y}_{\alpha,c}^l)^2] = \mathbb{E}_{\mathbf{x},\alpha}[(\mathbf{W}_{c,:}^l \mathbf{r}_f(\mathbf{x}^{l-1}, \boldsymbol{\alpha}))^2] \\
&= \sum_i (\tilde{\mathbf{W}}_{c,i}^l)^2 \lambda_i.
\end{aligned} \tag{47}$$

Given the assumption of standard initialization, biases are initialized with zeros and weights are initialized as $\mathbf{W}^l \stackrel{\theta^l}{\sim} \tilde{\mathbf{W}}^l \stackrel{\theta^l}{\sim} \mathcal{N}(0, \sqrt{2/R_l} \mathbf{I})$. Under A_{l-1} , it follows from Corollary 3 that $\text{Tr} \mathbf{G}[\mathbf{r}_f(\mathbf{x}^{l-1}, \boldsymbol{\alpha})] = R_l \nu_2(\mathbf{x}^{l-1}) > 0$ and thus $\text{Tr} \mathbf{G}[\mathbf{r}_f(\mathbf{x}^{l-1}, \boldsymbol{\alpha})] = \sum_i \lambda_i > 0$. Combined with Eq. (47), we get

$$\forall c: \mathbb{P}_{\theta^l|A_{l-1}}[\nu_{2,c}(\mathbf{y}^l) = 0] = 0. \tag{48}$$

Now we introduce the symmetric propagation from section B and we denote $A_{l,c} = \{\nu_{2,c}(\mathbf{x}^l) \neq 0\}$, $\tilde{A}_{l,c} = \{\nu_{2,c}(\tilde{\mathbf{x}}^l) \neq 0\}$ and $A'_{l,c} = \{\nu_{2,c}(\mathbf{x}^l) = 0\}$, $\tilde{A}'_{l,c} = \{\nu_{2,c}(\tilde{\mathbf{x}}^l) = 0\}$ the complementary events. It follows from Eq. (48) and Eq. (25) that $\mathbb{P}_{\theta^l|A_{l-1}}[A'_{l,c} \cap \tilde{A}'_{l,c}] = 0$. Furthermore by spherical symmetry of the propagation, $\nu_{2,c}(\mathbf{x}^l) \stackrel{\theta^l}{\sim} \nu_{2,c}(\tilde{\mathbf{x}}^l)$ and thus $\mathbb{P}_{\theta^l|A_{l-1}}[A'_{l,c}] = \mathbb{P}_{\theta^l|A_{l-1}}[\tilde{A}'_{l,c}]$. We

then get

$$\begin{aligned}\forall c : \mathbb{P}_{\theta^l|A_{l-1}}[A'_{l,c} \cup \tilde{A}'_{l+1,c}] &= \mathbb{P}_{\theta^l|A_{l-1}}[A'_{l,c}] + \mathbb{P}_{\theta^l|A_{l-1}}[\tilde{A}'_{l+1,c}] \leq 1, \\ \forall c : \mathbb{P}_{\theta^l|A_{l-1}}[A'_{l,c}] &\leq \frac{1}{2}.\end{aligned}$$

Since $A'_l = \bigcap_c A'_{l,c}$ and since the events $A'_{l,c}$ are independent conditionally on Θ^{l-1} and A_{l-1} , we conclude that $\mathbb{P}_{\theta^l|A_{l-1}}[A'_l] = \prod_c \mathbb{P}_{\theta^l|A_{l-1}}[A'_{l,c}] \leq 2^{-N_l}$.

The other side of the inequality is easier. If $\forall c, i_c : \mathbf{W}_{c,i_c}^l \leq 0$, it follows that $\forall \mathbf{x}, \boldsymbol{\alpha}, c : \mathbf{y}_{\boldsymbol{\alpha},c}^l < 0$, $\mathbf{x}_{\boldsymbol{\alpha},c}^l = (\mathbf{y}_{\boldsymbol{\alpha},c}^l)^+ = 0$, and thus $\nu_2(\mathbf{x}^l) = 0$. Therefore $\mathbb{P}_{\theta^l|A_{l-1}}[A'_l] \geq 2^{-K_l^d N_{l-1} N_l}$, since $K_l^d N_{l-1} N_l = R_l N_l$ is the number of elements in \mathbf{W}^l . We finally get

$$1 - 2^{-N_l} \leq \mathbb{P}_{\theta^l|A_{l-1}}[A_l] \leq 1 - 2^{-K_l^d N_{l-1} N_l}. \quad (49)$$

Since $\mathbb{P}[A_0] = 1$, due to the assumption $\nu_2(\mathbf{x}) = \mathbb{E}_{\mathbf{x}, \boldsymbol{\alpha}, c}[\mathbf{x}_{\boldsymbol{\alpha},c}^2] > 0$, it follows that $1 - 2^{-N_1} \leq \mathbb{P}[A_1] \leq 1 - 2^{-K_1^d N_0 N_1}$. This proves (i) for $l = 1$.

Now we proceed by induction and suppose that (i) is true for given $l - 1$. Using Eq. (49), we get

$$\prod_{k=1}^l (1 - 2^{-N_k}) \leq \mathbb{P}_{\Theta^l}[A_l] = \mathbb{P}_{\Theta^{l-1}}[A_{l-1}] \mathbb{P}_{\theta^l|A_{l-1}}[A_l] \leq \prod_{k=1}^l (1 - 2^{-K_k^d N_{k-1} N_k}),$$

meaning that (i) is true for l . (i) is thus true for all l . \square

Proof of (ii) for $\nu_2(\mathbf{x}^l)$. Again we denote $(\mathbf{e}_1, \dots, \mathbf{e}_{R_l})$ and $(\lambda_1, \dots, \lambda_{R_l})$ the eigenvectors and eigenvalues of $\mathbf{G}[\mathbf{r}_f(\mathbf{x}^{l-1}, \boldsymbol{\alpha})]$, and $\hat{\mathbf{W}}^l = \mathbf{W}^l(\mathbf{e}_1, \dots, \mathbf{e}_{R_l})$. We further define

$$u_c^l = \begin{cases} \frac{\nu_{2,c}(\mathbf{x}^l)}{\nu_{2,c}(\mathbf{x}^l) + \nu_{2,c}(\tilde{\mathbf{x}}^l)} & \text{if } \nu_{2,c}(\mathbf{x}^l) + \nu_{2,c}(\tilde{\mathbf{x}}^l) > 0 \\ 0 & \text{otherwise.} \end{cases}$$

Combining the definition of u_c^l with Eq. (25) and Eq. (47), we get

$$\begin{aligned}\forall c : \nu_{2,c}(\mathbf{x}^l) &= u_c^l \nu_{2,c}(\mathbf{y}^l), \\ \forall c : \nu_{2,c}(\mathbf{x}^l) &= u_c^l \sum_i (\hat{\mathbf{W}}_{c,i}^l)^2 \lambda_i = R_l \nu_2(\mathbf{x}^{l-1}) u_c^l \sum_i (\hat{\mathbf{W}}_{c,i}^l)^2 \hat{\lambda}_i,\end{aligned} \quad (50)$$

where we defined $\hat{\lambda}_i = \lambda_i / \sum_j \lambda_j$ and used $\sum_j \lambda_j = \text{Tr } \mathbf{G}[\mathbf{r}_f(\mathbf{x}^{l-1}, \boldsymbol{\alpha})] = R_l \nu_2(\mathbf{x}^{l-1})$ due to Corollary 3. The symmetric propagation gives

$$\begin{aligned}\forall c : \nu_{2,c}(\tilde{\mathbf{x}}^l) &= R_l \nu_2(\mathbf{x}^{l-1}) (1 - u_c^l) \sum_i (-\hat{\mathbf{W}}_{c,i}^l)^2 \hat{\lambda}_i, \\ \forall c : \nu_{2,c}(\mathbf{x}^l) + \nu_{2,c}(\tilde{\mathbf{x}}^l) &= R_l \nu_2(\mathbf{x}^{l-1}) \sum_i (\hat{\mathbf{W}}_{c,i}^l)^2 \hat{\lambda}_i.\end{aligned} \quad (51)$$

By symmetry of the propagation $\nu_{2,c}(\mathbf{x}^l) \stackrel{\theta^l}{\sim} \nu_{2,c}(\tilde{\mathbf{x}}^l)$. Combined with Eq. (51) and the assumption of standard initialization, we deduce

$$\begin{aligned}2\mathbb{E}_{\theta^l|A_{l-1}}[\nu_{2,c}(\mathbf{x}^l)] &= \mathbb{E}_{\theta^l|A_{l-1}}[\nu_{2,c}(\mathbf{x}^l) + \nu_{2,c}(\tilde{\mathbf{x}}^l)] \\ &= \mathbb{E}_{\theta^l|A_{l-1}}[R_l \nu_2(\mathbf{x}^{l-1}) \sum_i (\hat{\mathbf{W}}_{c,i}^l)^2 \hat{\lambda}_i] \\ &= R_l \nu_2(\mathbf{x}^{l-1}) \frac{2}{R_l} \sum_i \hat{\lambda}_i = 2\nu_2(\mathbf{x}^{l-1}).\end{aligned}$$

We then obtain $\forall c : \mathbb{E}_{\theta^l|A_{l-1}}[\nu_{2,c}(\mathbf{x}^l)] = \nu_2(\mathbf{x}^{l-1})$, and thus

$$\mathbb{E}_{\theta^l|A_{l-1}}[\nu_2(\mathbf{x}^l)] = \nu_2(\mathbf{x}^{l-1}), \quad \mathbb{E}_{\theta^l|A_{l-1}}[\delta \nu_2(\mathbf{x}^l)] = 1. \quad (52)$$

Both $\log x$ and $(\log x)^2$ are integrable at zero since $\int \log x dx = x \log x - x$ and $\int (\log x)^2 dx = x(\log x)^2 - 2x \log x + 2x$. Combined with Eq. (50), we deduce that $\log \delta \nu_2(\mathbf{x}^l)$ has well-defined conditional expectation and variance under A_l . Furthermore by Eq. (49), A_l has probability exponentially low in N_l , which gives

$$|\mathbb{E}_{\theta^l|A_l}[\delta \nu_2(\mathbf{x}^{l+1})] - 1| \ll 1, \quad \mathbb{E}_{\theta^l|A_l}[\log \delta \nu_2(\mathbf{x}^{l+1})] < 0, \quad (53)$$

where Eq. (53) is obtained by log-concavity. We now write the evolution of $\nu_2(\mathbf{x}^l)$ as

$$\log \nu_2(\mathbf{x}^l) - \nu_2(\mathbf{x}^0) = \sum_{k=1}^l \log \delta \nu_2(\mathbf{x}^k).$$

Let us define $X_k = \log \delta \nu_2(\mathbf{x}^k)$, $Y_k = \mathbb{E}_{\theta^k|A_l}[\log \delta \nu_2(\mathbf{x}^k)]$ and $Z_k = \log \delta \nu_2(\mathbf{x}^k) - \mathbb{E}_{\theta^k|A_l}[\log \delta \nu_2(\mathbf{x}^k)]$ and apply Lemma 10 for each l conditionally on A_l . Suppose there exist m_{\min} , m_{\max} , v_{\min} , v_{\max} , such that for each l we have conditionally on A_l that $\forall k \leq l$, $m_{\min} \leq -Y_k \leq m_{\max}$ and $v_{\min} \leq \text{Var}_{\theta^k|A_l}[Z_k] \leq v_{\max}$ which implies $v_{\min} \leq \text{Var}_{\theta^k|A_l}[Z_k] \leq v_{\max}$. Then we have $\forall k \leq l$, $-m_{\max} \leq Y_k \leq -m_{\min}$ and by Lemma 10 there exist sequences of random variables (m_l) and (s_l) such that $\forall l$ under A_l , s_l is centered and

$$\log \nu_2(\mathbf{x}^l) - \log \nu_2(\mathbf{x}^0) = l m_l + \sqrt{l} s_l, \quad -m_{\max} \leq m_l \leq -m_{\min}, \quad v_{\min} \leq \text{Var}_{\theta^l|A_l}[s_l] \leq v_{\max}.$$

By simply changing the variable m_l to $-m_l$, we get

$$\log \nu_2(\mathbf{x}^l) - \log \nu_2(\mathbf{x}^0) = -l m_l + \sqrt{l} s_l, \quad m_{\min} \leq m_l \leq m_{\max}, \quad v_{\min} \leq \text{Var}_{\theta^l|A_l}[s_l] \leq v_{\max}.$$

To obtain the bounds m_{\min} , m_{\max} , v_{\min} , v_{\max} , we consider extreme cases for u_c^l and $\sum_i (\hat{\mathbf{W}}_{c,i}^l)^2 \hat{\lambda}_i$ in Eq. (50). Denoting $\chi^2(N)$ the chi-square distribution with N degree of freedom, we obtain *minimum bounds* by considering $u_c^l \stackrel{\theta^l}{\sim} 1/2$ and $\sum_i (\hat{\mathbf{W}}_{c,i}^l)^2 \hat{\lambda}_i \stackrel{\theta^l}{\sim} \chi^2(R_l)$. This leads to $\log \delta \nu_2(\mathbf{x}^l) \stackrel{\theta^l}{\sim} \chi^2(N_l R_l)$. Denoting Bern(1/2) the Bernoulli distribution with $p = 1/2$, we obtain *maximum bounds* by considering $u_c \stackrel{\theta^l}{\sim} \text{Bern}(1/2)$ and $\sum_i (\hat{\mathbf{W}}_{c,i}^l)^2 \hat{\lambda}_i \stackrel{\theta^l}{\sim} \chi^2(1)$. The conditionality on A_l has highly negligible impact in practice.

Let us give an example in the fully-connected case with constant width $N_l = 100$. We then find numerically $m_{\min} \simeq 9.7 \times 10^{-5}$ and $v_{\min} \simeq 2.0 \times 10^{-4}$ as minimum bounds and $m_{\max} \simeq 2.5 \times 10^{-2}$ and $v_{\max} \simeq 5.2 \times 10^{-2}$ as maximum bounds. The length scale is experimentally close to $L_{\max} = 1/m_{\max} \simeq 40$ for $\nu_2(\mathbf{x}^l)$. \square

Proof of (ii) for $\mu_2(\mathbf{s}^l)$. Let us denote $B_l = \{\forall k \leq l, \forall c : \mu_{2,c}(\mathbf{s}^k) + \mu_{2,c}(\tilde{\mathbf{s}}^k) = \mu_{2,c}(\mathbf{t}^l)\}$ and $B_l' = \{\exists k \leq l, \exists c : \mu_{2,c}(\mathbf{s}^k) + \mu_{2,c}(\tilde{\mathbf{s}}^k) \neq \mu_{2,c}(\mathbf{t}^l)\} = \bigcup_{k=1}^l \{\exists c : \mu_{2,c}(\mathbf{s}^k) + \mu_{2,c}(\tilde{\mathbf{s}}^k) \neq \mu_{2,c}(\mathbf{t}^k)\}$ the complementary event. Eq. (28) gives $\forall k : \mathbb{P}_{\theta^k}[\{\exists c : \mu_{2,c}(\mathbf{s}^k) + \mu_{2,c}(\tilde{\mathbf{s}}^k) \neq \mu_{2,c}(\mathbf{t}^k)\}] = 0$. Since B_l' is the union of probability zero events, it follows that $\mathbb{P}_{\theta^l}(B_l') = 0$, $\mathbb{P}_{\theta^l}(B_l) = 1$, and thus $\mathbb{P}_{\theta^l}(B_l') = 0$, $\mathbb{P}_{\theta^l}(B_l) = 1$.

Since B_l has probability 1, the conditionality on B_l leaves moments of random variables unchanged. To see this, consider the moment of order p of a random variable x :

$$\begin{aligned} (\mathbb{E}_{\theta^l|A_l}[x^p] - \mathbb{E}_{\theta^l|A_l, B_l}[x^p])^2 &= \left(\mathbb{E}_{\theta^l|A_l}[x^p] - \frac{1}{\mathbb{P}_{\theta^l|A_l}(B_l)} \mathbb{E}_{\theta^l|A_l}[\mathbf{1}_{B_l} x^p] \right)^2 \\ &= \mathbb{E}_{\theta^l|A_l}[\mathbf{1}_{B_l'} x^p]^2 \\ &\leq \mathbb{E}_{\theta^l|A_l}[x^{2p}] \mathbb{P}_{\theta^l|A_l}[B_l'] = 0, \end{aligned} \quad (54)$$

where Eq. (54) is obtained with Cauchy-Schwarz inequality and the implicit assumption that x has well-defined moment of order $2p$. It follows that all arguments used in the proof of (i) remain valid, in particular regarding the distribution of \mathbf{W}^l . Conditionality on B_l , the proof then proceeds identically by simply replacing $\nu_2(\mathbf{x}^l)$ by $\mu_2(\mathbf{s}^l)$, \mathbf{y}^l by \mathbf{t}^l , \mathbf{G} by \mathbf{C} , and using the identity with μ_2 instead of ν_2 in Corollary 3. In particular, we have

$$\mathbb{E}_{\theta^l|A_{l-1},B_l}[\mu_2(\mathbf{s}^l)] = \mu_2(\mathbf{s}^{l-1}), \quad \mathbb{E}_{\theta^l|A_{l-1},B_l}[\delta\mu_2(\mathbf{s}^l)] = 1. \quad (55)$$

Furthermore under $A_l \cap B_l$, for the same positive constants $m_{\min}, m_{\max}, v_{\min}, v_{\max} > 0$ as previously defined, there exist sequences of random variables (m'_l) and (s'_l) such that $\forall l$ under $A_l \cap B_l$, s'_l is centered and

$$\log \mu_2(\mathbf{s}^l) = -lm'_l + \sqrt{l}s'_l, \quad m_{\min} \leq m'_l \leq m_{\max}, \quad v_{\min} \leq \text{Var}_{\theta^l|A_l,B_l}[s'_l] \leq v_{\max},$$

where we used the fact that $\mu_2(\mathbf{s}^0) = 1$ by the definition of \mathbf{s} . Now we extend m'_l and s'_l to B_l by setting m'_l to any value between m_{\min} and m_{\max} and s'_l such that $\log \mu_2(\mathbf{s}^l) = -lm'_l + \sqrt{l}s'_l$. The reasoning of Eq. (54) then ensures that $\mathbb{E}_{\theta^l|A_l,B_l}[s'_l] = \mathbb{E}_{\theta^l|A_l}[s'_l]$ and $\mathbb{E}_{\theta^l|A_l,B_l}[(s'_l)^2] = \mathbb{E}_{\theta^l|A_l}[(s'_l)^2]$, and $\text{Var}_{\theta^l|A_l,B_l}[s'_l] = \text{Var}_{\theta^l|A_l}[s'_l]$. It means that $\forall l$ under A_l , s'_l is centered and

$$\log \mu_2(\mathbf{s}^l) = -lm'_l + \sqrt{l}s'_l, \quad m_{\min} \leq m'_l \leq m_{\max}, \quad v_{\min} \leq \text{Var}_{\theta^l|A_l}[s'_l] \leq v_{\max}. \quad \square$$

E.3 RELATION TO THE TERMS \bar{m} , \underline{m} , \underline{c} DEFINED IN SECTION 4

Here we relate Theorem 1 to the terms \bar{m} , \underline{m} , \underline{c} defined in section 4. By Eq. (52), $|\mathbb{E}_{\theta^k|A_k}[\delta\nu_2(\mathbf{x}^k)] - 1| \ll 1$, and thus:

$$|\bar{m}[\nu_2(\mathbf{x}^k)]| = |\log \mathbb{E}_{\theta^k|A_k}[\delta\nu_2(\mathbf{x}^k)]| \simeq |\mathbb{E}_{\theta^k|A_k}[\delta\nu_2(\mathbf{x}^k)] - 1| \ll 1.$$

As in the proof of Theorem 1, we denote $B_l = \{\forall k \leq l, \forall c : \mu_{2,c}(\mathbf{s}^k) + \mu_{2,c}(\bar{\mathbf{s}}^k) = \mu_{2,c}(\mathbf{t}^l)\}$ and B'_l the complementary event. Then conditionally on A_k and B_k :

$$|\log \mathbb{E}_{\theta^k|A_k,B_k}[\delta\nu_2(\mathbf{x}^k)]| \simeq |\mathbb{E}_{\theta^k|A_k,B_k}[\delta\nu_2(\mathbf{x}^k)] - 1| \ll 1.$$

The reasoning of Eq. (54) can be applied to $\delta\nu_2(\mathbf{x}^k)$, which results in $\mathbb{E}_{\theta^k|A_k}[\delta\nu_2(\mathbf{x}^k)] = \mathbb{E}_{\theta^k|A_k,B_k}[\delta\nu_2(\mathbf{x}^k)]$. Therefore we also get

$$|\bar{m}[\mu_2(\mathbf{s}^k)]| = |\log \mathbb{E}_{\theta^k|A_k}[\delta\mu_2(\mathbf{s}^k)]| \ll 1.$$

The terms $\bar{m}[\nu_2(\mathbf{x}^k)]$ and $\bar{m}[\mu_2(\mathbf{s}^k)]$ are thus vanishing and the evolution is dominated by the terms $\bar{m}[\nu_2(\mathbf{x}^k)] < 0$, $\bar{m}[\mu_2(\mathbf{s}^k)] < 0$. These terms correspond to Y_k in the proof of Theorem 1 (ii).

E.4 CONVERGENCE IN PROBABILITY TO ZERO

Corollary 11. *Conditionally on A_l the variables $\nu_2(\mathbf{x}^l)$ and $\mu_2(\mathbf{s}^l)$ still converge in probability to zero:*

$$\forall \epsilon : \mathbb{P}_{\theta^l|A_l}[|\nu_2(\mathbf{x}^l)| > \epsilon] \rightarrow 0, \quad \forall \epsilon : \mathbb{P}_{\theta^l|A_l}[|\mu_2(\mathbf{s}^l)| > \epsilon] \rightarrow 0.$$

Proof. Consider a given ϵ and the evolution of $\nu_2(\mathbf{x}^l)$. From theorem 1 (ii), we can write under A_l : $\log \nu_2(\mathbf{x}^l) = -lm_l + \sqrt{l}s_l + \log \nu_2(\mathbf{x}^0)$, with s_l centered and $m_{\min} \leq m_l \leq m_{\max}$, $v_{\min} \leq \text{Var}_{\theta^l|A_l}[s_l] \leq v_{\max}$. Therefore:

$$\begin{aligned} \mathbb{P}_{\theta^l|A_l}[|\nu_2(\mathbf{x}^l)| > \epsilon] &= \mathbb{P}_{\theta^l|A_l}[\log \nu_2(\mathbf{x}^l) > \log \epsilon] = \mathbb{P}_{\theta^l|A_l}[\sqrt{l}s_l > lm_l + \log \epsilon - \log \nu_2(\mathbf{x}^0)] \\ &\leq \mathbb{P}_{\theta^l|A_l}\left[s_l > \frac{1}{\sqrt{l}}(lm_l + \log \epsilon - \log \nu_2(\mathbf{x}^0))\right] \\ &\leq \mathbb{P}_{\theta^l|A_l}\left[|s_l| > \frac{1}{\sqrt{l}}(lm_{\min} + \log \epsilon - \log \nu_2(\mathbf{x}^0))\right]. \end{aligned} \quad (56)$$

Chebyshev's inequality on the centered random variable s_l then gives:

$$\begin{aligned} \mathbb{P}_{\Theta^l|A_l} [|\nu_2(\mathbf{x}^l)| > \epsilon] &\leq \frac{l}{(lm_{\min} + \log \epsilon - \log \nu_2(\mathbf{x}^0))^2} \text{Var}_{\Theta^l|A_l} [s_l] \\ &\leq \frac{l}{(lm_{\min} + \log \epsilon - \log \nu_2(\mathbf{x}^0))^2} v_{\max} \sim \frac{1}{l} \frac{v_{\max}}{m_{\min}^2} \end{aligned}$$

where \sim denotes the equivalence for large l . It follows that $\mathbb{P}_{\Theta^l|A_l} [|\nu_2(\mathbf{x}^l)| > \epsilon] \rightarrow 0$. The same analysis applied to $\mu_2(\mathbf{s}^l)$ gives $\mathbb{P}_{\Theta^l|A_l} [|\mu_2(\mathbf{s}^l)| > \epsilon] \rightarrow 0$. \square

F NORMALIZED SENSITIVITY INCREMENTS OF VANILLA NETWORKS

F.1 PROOF OF THEOREM 2

Theorem 2. Normalized Sensitivity increments of vanilla networks. *Under A_{l-1} , the dominating term in the evolution of the normalized sensitivity is:*

$$\delta\zeta^l \simeq \exp\left(\overline{m}_{\text{vanilla}}[\zeta^l]\right) = \left(1 - \mathbb{E}_{\mathbf{c}, \theta^l|A_{l-1}} \left[\frac{\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-})}{\mu_2(\mathbf{x}^{l-1})} \right]\right)^{-1/2}, \quad (57)$$

where $\mathbf{y}^{l,+} = \max(\mathbf{y}^l, 0)$ and $\mathbf{y}^{l,-} = \max(-\mathbf{y}^l, 0)$.

Proof. The dominating term in the evolution of ζ^l is $\frac{1}{2}(\overline{m}[\mu_2(\mathbf{s}^l)] - \overline{m}[\mu_2(\mathbf{x}^l)])$. The terms $\overline{m}[\mu_2(\mathbf{s}^l)]$ and $\overline{m}[\mu_2(\mathbf{x}^l)]$ are simply obtained by considering $\mathbb{E}_{\theta^l}[\delta\mu_2(\mathbf{s}^l)]$ and $\mathbb{E}_{\theta^l}[\delta\mu_2(\mathbf{x}^l)]$.

By Eq. (53) in the proof of Theorem 1: $\mathbb{E}_{\theta^l|A_{l-1}, B_l}[\delta\mu_2(\mathbf{s}^l)] = 1$. By replicating the reasoning of Eq. (54), this further gives $\mathbb{E}_{\theta^l|A_{l-1}}[\delta\mu_2(\mathbf{s}^l)] = 1$, and thus $\overline{m}[\mu_2(\mathbf{s}^l)] = \log \mathbb{E}_{\theta^l|A_{l-1}}[\delta\mu_2(\mathbf{s}^l)] = 0$.

Next we turn to the term $\overline{m}[\mu_2(\mathbf{x}^l)]$. Again we use the definitions and notations from section B. We further denote $(\mathbf{e}_1, \dots, \mathbf{e}_{R_l})$ and $(\lambda_1, \dots, \lambda_{R_l})$ respectively the orthogonal eigenvectors and eigenvalues of $\mathbf{C}[\mathbf{r}_f(\mathbf{x}^{l-1}, \boldsymbol{\alpha})]$ and $\hat{\mathbf{W}}^l = \mathbf{W}^l(\mathbf{e}_1, \dots, \mathbf{e}_{R_l})$. Using these notations, we get

$$\begin{aligned} \forall \mathbf{c}: \mu_{2,c}(\mathbf{y}^l) &= \mathbb{E}_{\mathbf{x}, \boldsymbol{\alpha}} \left[\hat{\mathbf{f}}_m(\mathbf{y}^l, \boldsymbol{\alpha})_c^2 \right] = \mathbb{E}_{\mathbf{x}, \boldsymbol{\alpha}} \left[(\mathbf{W}_{c,:}^l \hat{\mathbf{r}}_f(\mathbf{x}^{l-1}, \boldsymbol{\alpha}))^2 \right] \\ &= \sum_i (\hat{\mathbf{W}}_{c,i}^l)^2 \lambda_i. \end{aligned} \quad (58)$$

Then due to $\mathbf{W}^l \stackrel{\theta^l}{\sim} \hat{\mathbf{W}}^l \stackrel{\theta^l}{\sim} \mathcal{N}(0, \sqrt{2/R_l} \mathbf{I})$;

$$\begin{aligned} \mathbb{E}_{\theta^l|A_{l-1}}[\mu_{2,c}(\mathbf{y}^l)] &= \frac{2}{R_l} \sum_i \lambda_i = \frac{2}{R_l} \text{Tr} \mathbf{C}[\mathbf{r}_f(\mathbf{x}^{l-1}, \boldsymbol{\alpha})] \\ &= 2\mu_2(\mathbf{x}^{l-1}). \end{aligned} \quad (59)$$

where Eq. (59) follows from Corollary 3. Furthermore the symmetric propagation gives:

$$\begin{aligned} \mu_{2,c}(\mathbf{x}^l) + \mu_{2,c}(\tilde{\mathbf{x}}^l) &= \mathbb{E}_{\mathbf{x}, \boldsymbol{\alpha}} [(\mathbf{y}_{\alpha,c}^{l,+})^2] - \mathbb{E}_{\mathbf{x}, \boldsymbol{\alpha}} [\mathbf{y}_{\alpha,c}^{l,+}]^2 + \mathbb{E}_{\mathbf{x}, \boldsymbol{\alpha}} [(\mathbf{y}_{\alpha,c}^{l,-})^2] - \mathbb{E}_{\mathbf{x}, \boldsymbol{\alpha}} [\mathbf{y}_{\alpha,c}^{l,-}]^2 \\ &= \nu_{2,c}(\mathbf{y}^{l,+}) - \nu_{1,c}(\mathbf{y}^{l,+})^2 + \nu_{2,c}(\mathbf{y}^{l,-}) - \nu_{1,c}(\mathbf{y}^{l,-})^2 \\ &= \nu_{2,c}(\mathbf{y}^l) - \left(\nu_{1,c}(\mathbf{y}^{l,+})^2 + \nu_{1,c}(\mathbf{y}^{l,-})^2 \right) \end{aligned} \quad (60)$$

We have $\nu_{1,c}(\mathbf{y}^l) = \nu_{1,c}(\mathbf{y}^{l,+}) - \nu_{1,c}(\mathbf{y}^{l,-})$ and thus $\nu_{1,c}(\mathbf{y}^l)^2 = \nu_{1,c}(\mathbf{y}^{l,+})^2 + \nu_{1,c}(\mathbf{y}^{l,-})^2 - 2\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-})$. We can then rewrite Eq. (60) as

$$\begin{aligned} \mu_{2,c}(\mathbf{x}^l) + \mu_{2,c}(\tilde{\mathbf{x}}^l) &= \nu_{2,c}(\mathbf{y}^l) - \nu_{1,c}(\mathbf{y}^l)^2 - 2\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-}) \\ &= \mu_{2,c}(\mathbf{y}^l) - 2\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-}) \end{aligned} \quad (61)$$

Combining Eq. (59) and Eq. (61):

$$\begin{aligned}\mathbb{E}_{\theta^l|A_{l-1}}[\mu_{2,c}(\mathbf{x}^l) + \mu_{2,c}(\tilde{\mathbf{x}}^l)] &= 2\mu_2(\mathbf{x}^{l-1}) - 2\mathbb{E}_{\theta^l|A_{l-1}}[\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-})], \\ 2\mathbb{E}_{\theta^l|A_{l-1}}[\mu_{2,c}(\mathbf{x}^l)] &= 2\mu_2(\mathbf{x}^{l-1}) - 2\mathbb{E}_{\theta^l|A_{l-1}}[\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-})], \\ \mathbb{E}_{\theta^l|A_{l-1}}[\mu_{2,c}(\mathbf{x}^l)] &= \mu_2(\mathbf{x}^{l-1}) \left(1 - \mathbb{E}_{\theta^l|A_{l-1}} \left[\frac{\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-})}{\mu_2(\mathbf{x}^{l-1})} \right] \right).\end{aligned}\quad (62)$$

where Eq. (62) is obtained by spherical symmetry of the propagation. We finally get

$$\begin{aligned}\mathbb{E}_{\theta^l|A_{l-1}}[\mu_2(\mathbf{x}^l)] &= \mathbb{E}_{\theta^l|A_{l-1}}[\mathbb{E}_c[\mu_{2,c}(\mathbf{x}^l)]] = \mathbb{E}_c[\mathbb{E}_{\theta^l|A_{l-1}}[\mu_{2,c}(\mathbf{x}^l)]] \\ &= \mu_2(\mathbf{x}^{l-1}) \left(1 - \mathbb{E}_{c,\theta^l|A_{l-1}} \left[\frac{\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-})}{\mu_2(\mathbf{x}^{l-1})} \right] \right), \\ \mathbb{E}_{\theta^l|A_{l-1}}[\delta\mu_2(\mathbf{x}^{l-1})] &= 1 - \mathbb{E}_{c,\theta^l|A_{l-1}} \left[\frac{\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-})}{\mu_2(\mathbf{x}^{l-1})} \right].\end{aligned}$$

Combined with $\mathbb{E}_{\theta^l|A_{l-1}}[\delta\mu_2(\mathbf{s}^l)] = 1$,

$$\begin{aligned}\delta\zeta^l &\simeq \exp\left(\overline{m}_{\text{vanilla}}[\zeta^l]\right) = \left(\frac{\mathbb{E}_{\theta^l|A_{l-1}}[\delta\mu_2(\mathbf{s}^l)]}{\mathbb{E}_{\theta^l|A_{l-1}}[\delta\mu_2(\mathbf{x}^l)]} \right)^{1/2} \\ &= \left(1 - \mathbb{E}_{c,\theta^l|A_{l-1}} \left[\frac{\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-})}{\mu_2(\mathbf{x}^{l-1})} \right] \right)^{-1/2}.\end{aligned}$$

□

F.2 IF ζ^l HAS DRIFT LARGER THAN DIFFUSION, THEN $\mu_2(\mathbf{x}^l) / \nu_2(\mathbf{x}^l)$ CONVERGES IN PROBABILITY TO ZERO

We only need to slightly adapt the reasoning of section E.4. From theorem 1, we can write under A_l : $\log \nu_2(\mathbf{x}^l) = -lm_l + \sqrt{l}s_l + \log \nu_2(\mathbf{x}^0)$, with s_l centered and $m_{\min} \leq m_l \leq m_{\max}$, $v_{\min} \leq \text{Var}_{\theta^l|A_l}[s_l] \leq v_{\max}$. We can also write $\log \mu_2(\mathbf{s}^l) = -lm'_l + \sqrt{l}s'_l$, with s'_l centered and $m_{\min} \leq m'_l \leq m_{\max}$ and $v_{\min} \leq \text{Var}_{\theta^l|A_l}[s'_l] \leq v_{\max}$. We further suppose that there is an event D with $\mathbb{P}(D) > 0$ under which ζ^l has *drift larger than diffusion*. To make it precise, this means that $\exists m > \frac{1}{2}(m_{\max} - m_{\min})$ such that for l large enough: $\log \delta\zeta^l \geq m$, and thus $\exists c \in \mathbb{R}$ such that for $\forall l$: $\log \zeta^l \geq c + lm$.

The ratio $\mu_2(\mathbf{x}^l) / \nu_2(\mathbf{x}^l)$ can be expressed as

$$\frac{\mu_2(\mathbf{x}^l)}{\nu_2(\mathbf{x}^l)} = \frac{\mu_2(\mathbf{x}^l)}{\mu_2(\mathbf{s}^l)\mu_2(\mathbf{x}^0)} \frac{\mu_2(\mathbf{s}^l)\mu_2(\mathbf{x}^0)}{\nu_2(\mathbf{x}^l)} = \frac{1}{(\zeta^l)^2} \frac{\mu_2(\mathbf{s}^l)\mu_2(\mathbf{x}^0)}{\nu_2(\mathbf{x}^l)},$$

which gives with logarithms,

$$\begin{aligned}\log \mu_2(\mathbf{x}^l) - \log \nu_2(\mathbf{x}^l) &= -2 \log \zeta^l + \log \mu_2(\mathbf{s}^l) - \log \nu_2(\mathbf{x}^l) + \log \mu_2(\mathbf{x}^0) \\ &\leq -2c - 2lm - lm_{\min} + \sqrt{l}s'_l + lm_{\max} - \sqrt{l}s_l - \log \nu_2(\mathbf{x}^0) + \log \mu_2(\mathbf{x}^0) \\ &\leq C - lM + \sqrt{l}s_l,\end{aligned}$$

where we denoted $C = -2c - \log \nu_2(\mathbf{x}^0) + \log \mu_2(\mathbf{x}^0)$, $M = 2m + m_{\min} - m_{\max} > 0$ and $s_l = s'_l - s_l$. The variance of s_l is bounded as

$$\begin{aligned}\mathbb{E}_{\theta^l|A_l}[s_l^2] &= \text{Var}_{\theta^l|A_l}[s_l] + \text{Var}_{\theta^l|A_l}[s'_l] - 2\mathbb{E}_{\theta^l|A_l}[s_l s'_l] \\ &\leq \text{Var}_{\theta^l|A_l}[s_l] + \text{Var}_{\theta^l|A_l}[s'_l] + 2\text{Var}_{\theta^l|A_l}[s_l]^{1/2} \text{Var}_{\theta^l|A_l}[s'_l]^{1/2} \leq 4v_{\max} \\ \text{Var}_{\theta^l|A_l,D}[s_l] &= \mathbb{E}_{\theta^l|A_l,D}[s_l^2] = \frac{1}{\mathbb{P}(D)} \mathbb{E}_{\theta^l|A_l}[\mathbf{1}_D s_l^2] \leq \frac{1}{\mathbb{P}(D)} 4v_{\max}.\end{aligned}$$

Now for given ϵ :

$$\begin{aligned} \mathbb{P}_{\Theta^l|A_l,D} \left[\frac{\mu_2(\mathbf{x}^l)}{\nu_2(\mathbf{x}^l)} > \epsilon \right] &= \mathbb{P}_{\Theta^l|A_l,D} [\log \mu_2(\mathbf{x}^l) - \log \nu_2(\mathbf{x}^l) > \log \epsilon] \\ &\leq \mathbb{P}_{\Theta^l|A_l,D} \left[s_l > \frac{1}{\sqrt{l}} (lM + \log \epsilon - C) \right] \end{aligned}$$

Chebyshev's inequality on the centered random variable s_l further gives

$$\begin{aligned} \mathbb{P}_{\Theta^l|A_l,D} \left[\frac{\mu_2(\mathbf{x}^l)}{\nu_2(\mathbf{x}^l)} > \epsilon \right] &\leq \frac{l}{(lM + \log \epsilon - C)^2} \text{Var}_{\Theta^l|A_l,D}[s_l] \\ &\leq \frac{l}{(lM + \log \epsilon - C)^2} \frac{1}{\mathbb{P}(D)} 4v_{\max} \sim \frac{1}{l} \frac{4v_{\max}}{M^2 \mathbb{P}(D)}, \end{aligned}$$

proving that under D the ratio $\mu_2(\mathbf{x}^l) / \nu_2(\mathbf{x}^l)$ converges in probability to zero.

F.3 LIMITS OF SIGNAL DISTRIBUTION

Proposition 12. *Suppose that*

$$\mathbb{P}_{c,\theta^l} \left[\left\{ \min(\nu_{1,c}(\mathbf{y}^{l,+}), \nu_{1,c}(\mathbf{y}^{l,-})) = 0 \right\} \cap \left\{ \max(\nu_{1,c}(\mathbf{y}^{l,+}), \nu_{1,c}(\mathbf{y}^{l,-})) > 0 \right\} \right] = 1.$$

Then $\mathbf{f}_m(\mathbf{y}^l, \boldsymbol{\alpha})$ concentrates on the semi-line generated by its average vector $(\nu_{1,c}(\mathbf{y}^l))_{1 \leq c \leq N_l}$.

Proof. Let us consider the feature map vectors $\mathbf{f}_m(\mathbf{x}^{l-1}, \boldsymbol{\alpha})$ and receptive field vectors $\mathbf{r}_f(\mathbf{x}^{l-1}, \boldsymbol{\alpha})$. Due to the statistics-preserving property of Corollary 3, for each channel c and index $i_c \in \mathcal{I}_c^l$, $\mathbf{x}_{\alpha,c}^{l-1} \stackrel{\mathbf{x}, \boldsymbol{\alpha}}{\sim} RF(\mathbf{x}^{l-1})_{\alpha, i_c}$ and thus $\mathbb{E}_{\mathbf{x}, \boldsymbol{\alpha}}[\mathbf{f}_m(\mathbf{x}^{l-1}, \boldsymbol{\alpha})_c] = \mathbb{E}_{\mathbf{x}, \boldsymbol{\alpha}}[\mathbf{r}_f(\mathbf{x}^{l-1}, \boldsymbol{\alpha})_{i_c}]$.

Now let us reason by contradiction and suppose that there exists a direction \mathbf{e} which is orthogonal to $(\mathbb{E}_{\mathbf{x}, \boldsymbol{\alpha}}[\mathbf{r}_f(\mathbf{x}^{l-1}, \boldsymbol{\alpha})_{i_c}])_{1 \leq i_c \leq R_l}$ and with non-zero variance $v > 0$. Consider the weight matrix \mathbf{W}^l which projects on direction \mathbf{e} for given channel c . For this direction, we have $\nu_{1,c}(\mathbf{y}^l) = 0$, $\nu_{2,c}(\mathbf{y}^l) > 0$. In turn, this implies $\nu_{1,c}(\mathbf{y}^{l,+}) = \nu_{1,c}(\mathbf{y}^{l,-})$ and $\nu_{1,c}(\mathbf{y}^{l,+}) + \nu_{1,c}(\mathbf{y}^{l,-}) > 0$, which further gives $\min(\nu_{1,c}(\mathbf{y}^{l,+}), \nu_{1,c}(\mathbf{y}^{l,-})) > 0$. By continuity, $\min(\nu_{1,c}(\mathbf{y}^{l,+}), \nu_{1,c}(\mathbf{y}^{l,-})) > 0$ in a small neighborhood for the sampling of the weights \mathbf{W}^l , which contradicts the hypothesis. It follows that $\mathbf{r}_f(\mathbf{x}^l, \boldsymbol{\alpha})$ concentrates on the direction generated by $(\mathbb{E}_{\mathbf{x}, \boldsymbol{\alpha}}[\mathbf{r}_f(\mathbf{x}^{l-1}, \boldsymbol{\alpha})_{i_c}])_{1 \leq i_c \leq R_l}$.

Now consider the weight matrix \mathbf{W}^l which projects on the direction generated by $(\mathbb{E}_{\mathbf{x}, \boldsymbol{\alpha}}[\mathbf{r}_f(\mathbf{x}^{l-1}, \boldsymbol{\alpha})_{i_c}])_{1 \leq i_c \leq R_l}$ for given channel c . A similar argument gives $\min(\nu_{1,c}(\mathbf{y}^{l,+}), \nu_{1,c}(\mathbf{y}^{l,-})) = 0$, and thus that $\mathbf{f}_m(\mathbf{y}^l, \boldsymbol{\alpha})_c$ either concentrates in \mathbb{R}^+ or concentrates in \mathbb{R}^- . It follows that $\mathbf{r}_f(\mathbf{x}^{l-1}, \boldsymbol{\alpha})$ concentrates on the semi-line generated by $(\mathbb{E}_{\mathbf{x}, \boldsymbol{\alpha}}[\mathbf{r}_f(\mathbf{x}^{l-1}, \boldsymbol{\alpha})_{i_c}])_{1 \leq i_c \leq R_l}$. Under standard initialization, the image $\mathbf{f}_m(\mathbf{y}^l, \boldsymbol{\alpha})$ of $\mathbf{r}_f(\mathbf{x}^{l-1}, \boldsymbol{\alpha})$ by the affine transform $\mathbf{f}_m(\mathbf{y}^l, \boldsymbol{\alpha}) = \mathbf{W}^l \mathbf{r}_f(\mathbf{x}^{l-1}, \boldsymbol{\alpha}) + \mathbf{b}^l = \mathbf{W}^l \mathbf{r}_f(\mathbf{x}^{l-1}, \boldsymbol{\alpha})$ thus concentrates on the semi-line generated by its average vector $(\nu_{1,c}(\mathbf{y}^l))_{1 \leq c \leq N_l}$. \square

G NORMALIZED SENSITIVITY INCREMENTS OF BATCH-NORMALIZED FEEDFORWARD NETS

G.1 PROOF OF THEOREM 3

Theorem 3. *Normalized Sensitivity increments of batch-normalized feedforward nets. The dominating term in the evolution of ζ^l can be decomposed as the sum of a term $\overline{m}_{BN}[\zeta^l]$ due to batch normalization and a term $\overline{m}_\phi[\zeta^l]$ due to the nonlinearity ϕ :*

$$\begin{aligned}
\exp\left(\overline{m}_{BN}[\zeta^l]\right) &= \left(\frac{\mu_2(\mathbf{s}^{l-1})}{\mu_2(\mathbf{x}^{l-1})}\right)^{-1/2} \mathbb{E}_{\mathbf{c},\theta^l} \left[\frac{\mu_{2,c}(\mathbf{t}^l)}{\mu_{2,c}(\mathbf{y}^l)}\right]^{1/2}, \\
\exp\left(\overline{m}_\phi[\zeta^l]\right) &= \left(1 - 2\mathbb{E}_{\mathbf{c},\theta^l} \left[\nu_{1,c}(\mathbf{z}^{l,+})\nu_{1,c}(\mathbf{z}^{l,-})\right]\right)^{-1/2}, \\
\delta\zeta^l &\simeq \exp\left(\overline{m}_{BN/FF}[\zeta^l]\right) = \exp\left(\overline{m}_{BN}[\zeta^l] + \overline{m}_\phi[\zeta^l]\right).
\end{aligned}$$

Proof. First let us decompose the dominating term as the product of two terms:

$$\begin{aligned}
\exp\left(\overline{m}_{BN}[\zeta^l]\right) &= \left(\frac{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{u}^l)]}{\mu_2(\mathbf{s}^{l-1})}\right)^{1/2} \left(\frac{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{z}^l)]}{\mu_2(\mathbf{x}^{l-1})}\right)^{-1/2}, \\
\exp\left(\overline{m}_\phi[\zeta^l]\right) &= \left(\frac{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{s}^l)]}{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{u}^l)]}\right)^{1/2} \left(\frac{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{x}^l)]}{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{z}^l)]}\right)^{-1/2}, \\
\exp\left(\overline{m}_{BN/FF}[\zeta^l]\right) &= \left(\frac{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{s}^l)]}{\mu_2(\mathbf{s}^{l-1})}\right)^{1/2} \left(\frac{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{x}^l)]}{\mu_2(\mathbf{x}^{l-1})}\right)^{-1/2} \\
&= \exp\left(\overline{m}_{BN}[\zeta^l]\right) \exp\left(\overline{m}_\phi[\zeta^l]\right).
\end{aligned}$$

$\overline{m}_{BN}[\zeta^l]$ is a dominating term in the evolution of ζ^l from $(\mathbf{x}^{l-1}, \mathbf{s}^{l-1})$ to $(\mathbf{z}^l, \mathbf{u}^l)$, while $\overline{m}_\phi[\zeta^l]$ is a dominating term in the evolution of ζ^l from $(\mathbf{z}^l, \mathbf{u}^l)$ to $(\mathbf{x}^l, \mathbf{s}^l)$. These terms can be seen as the contribution to $\overline{m}_{BN/FF}[\zeta^l]$ of respectively batch normalization and ϕ . Now let us explicitate both terms.

Term $\exp\left(\overline{m}_{BN}[\zeta^l]\right)$. First we note that batch normalization directly gives $\mu_2(\mathbf{z}^l) = 1$ and thus $\mathbb{E}_{\theta^l}[\mu_2(\mathbf{z}^l)] = 1$. Now let us explicitate $\mathbb{E}_{\theta^l}[\mu_2(\mathbf{u}^l)]$:

$$\begin{aligned}
\forall \mathbf{c} : \mathbf{u}^l_{:,c} &= \frac{\mathbf{t}^l_{:,c}}{\mu_{2,c}(\mathbf{y}^l)^{1/2}}, \quad \forall \mathbf{c} : \mu_{2,c}(\mathbf{u}^l) = \frac{\mu_{2,c}(\mathbf{t}^l)}{\mu_{2,c}(\mathbf{y}^l)}, \\
\mathbb{E}_{\theta^l}[\mu_2(\mathbf{u}^l)] &= \mathbb{E}_{\mathbf{c},\theta^l}[\mu_{2,c}(\mathbf{u}^l)] = \mathbb{E}_{\mathbf{c},\theta^l} \left[\frac{\mu_{2,c}(\mathbf{t}^l)}{\mu_{2,c}(\mathbf{y}^l)}\right]
\end{aligned}$$

All together, we get

$$\begin{aligned}
\exp\left(\overline{m}_{BN}[\zeta^l]\right) &= \left(\frac{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{u}^l)]}{\mu_2(\mathbf{s}^{l-1})}\right)^{1/2} \left(\frac{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{z}^l)]}{\mu_2(\mathbf{x}^{l-1})}\right)^{-1/2} \\
&= \left(\frac{\mu_2(\mathbf{s}^{l-1})}{\mu_2(\mathbf{x}^{l-1})}\right)^{-1/2} \mathbb{E}_{\mathbf{c},\theta^l} \left[\frac{\mu_{2,c}(\mathbf{t}^l)}{\mu_{2,c}(\mathbf{y}^l)}\right]^{1/2}
\end{aligned}$$

Term $\exp\left(\overline{m}_\phi[\zeta^l]\right)$. We consider the symmetric propagation for batch-normalized feedforward nets and again we denote $B_l = \{\forall k \leq l, \forall \mathbf{c} : \mu_{2,c}(\mathbf{s}^k) + \mu_{2,c}(\tilde{\mathbf{s}}^k) = \mu_{2,c}(\mathbf{u}^k)\} = \bigcap_{k=1}^l \{\forall \mathbf{c} : \mu_{2,c}(\mathbf{s}^k) + \mu_{2,c}(\tilde{\mathbf{s}}^k) = \mu_{2,c}(\mathbf{u}^k)\}$ and B'_l the complementary event. By Eq. (33): $\mathbb{P}_{\Theta^l}(B_l) = 1$, $\mathbb{P}_{\Theta^l}(B'_l) = 0$, and replicating the reasoning of Eq. (54) we deduce

$$\begin{aligned}
\mathbb{E}_{\theta^l}[\mu_2(\mathbf{s}^l)] + \mathbb{E}_{\theta^l}[\mu_2(\tilde{\mathbf{s}}^l)] &= \mathbb{E}_{\theta^l|B_l}[\mu_2(\mathbf{s}^l)] + \mathbb{E}_{\theta^l|B_l}[\mu_2(\tilde{\mathbf{s}}^l)] = \mathbb{E}_{\theta^l|B_l}[\mu_2(\mathbf{u}^l)] = \mathbb{E}_{\theta^l}[\mu_2(\mathbf{u}^l)], \\
2\mathbb{E}_{\theta^l}[\mu_2(\mathbf{s}^l)] &= \mathbb{E}_{\theta^l}[\mu_2(\mathbf{u}^l)], \tag{63}
\end{aligned}$$

where Eq. (63) follows from spherical symmetry of the propagation. Now we turn to the symmetric propagation of the signal:

$$\begin{aligned}
\mu_{2,c}(\mathbf{x}^l) + \mu_{2,c}(\tilde{\mathbf{x}}^l) &= \mathbb{E}_{\mathbf{x},\alpha}[(\mathbf{z}^{l,+}_{\alpha,c})^2] - \mathbb{E}_{\mathbf{x},\alpha}[\mathbf{z}^{l,+}_{\alpha,c}]^2 + \mathbb{E}_{\mathbf{x},\alpha}[(\mathbf{z}^{l,-}_{\alpha,c})^2] - \mathbb{E}_{\mathbf{x},\alpha}[\mathbf{z}^{l,-}_{\alpha,c}]^2. \tag{64} \\
&= \nu_{2,c}(\mathbf{z}^{l,+}) - \nu_{1,c}(\mathbf{z}^{l,+})^2 + \nu_{2,c}(\mathbf{z}^{l,-}) - \nu_{1,c}(\mathbf{z}^{l,-})^2 \\
&= \nu_{2,c}(\mathbf{z}^l) - \left(\nu_{1,c}(\mathbf{z}^{l,+})^2 + \nu_{1,c}(\mathbf{z}^{l,-})^2\right),
\end{aligned}$$

where Eq. (64) follows from Eq. (31). Due to the constraints imposed by batch normalization, $\nu_{1,c}(\mathbf{z}^l) = 0$ and $\nu_{2,c}(\mathbf{z}^l) = 1$, it follows that

$$\mu_{2,c}(\mathbf{x}^l) + \mu_{2,c}(\tilde{\mathbf{x}}^l) = 1 - \left(\nu_{1,c}(\mathbf{z}^{l,+})^2 + \nu_{1,c}(\mathbf{z}^{l,-})^2 \right). \quad (65)$$

$$\begin{aligned} \nu_{1,c}(\mathbf{z}^l) &= \nu_{1,c}(\mathbf{z}^{l,+}) - \nu_{1,c}(\mathbf{z}^{l,-}) = 0, \\ \left(\nu_{1,c}(\mathbf{z}^{l,+}) - \nu_{1,c}(\mathbf{z}^{l,-}) \right)^2 &= \nu_{1,c}(\mathbf{z}^{l,+})^2 + \nu_{1,c}(\mathbf{z}^{l,-})^2 - 2\nu_{1,c}(\mathbf{z}^{l,+})\nu_{1,c}(\mathbf{z}^{l,-}) = 0. \end{aligned} \quad (66)$$

Using Eq. (65), Eq. (66) and the symmetry of the propagation,

$$\begin{aligned} \mu_{2,c}(\mathbf{x}^l) + \mu_{2,c}(\tilde{\mathbf{x}}^l) &= 1 - 2\nu_{1,c}(\mathbf{z}^{l,+})\nu_{1,c}(\mathbf{z}^{l,-}), \\ 2\mathbb{E}_{\theta^l}[\mu_{2,c}(\mathbf{x}^l)] &= 1 - 2\mathbb{E}_{c,\theta^l} \left[\nu_{1,c}(\mathbf{z}^{l,+})\nu_{1,c}(\mathbf{z}^{l,-}) \right]. \end{aligned} \quad (67)$$

We finally combine Eq. (63) and Eq. (67):

$$\begin{aligned} \exp(\bar{m}_\phi[S^l]) &= \left(\frac{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{s}^l)]}{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{u}^l)]} \right)^{1/2} \left(\frac{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{x}^l)]}{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{z}^l)]} \right)^{-1/2}, \\ &= \left(1 - 2\mathbb{E}_{c,\theta^l} \left[\nu_{1,c}(\mathbf{z}^{l,+})\nu_{1,c}(\mathbf{z}^{l,-}) \right] \right)^{-1/2}. \end{aligned}$$

□

G.2 $\exp(\bar{m}_{BN}[\zeta^1]) > 1$ IN THE FIRST STEP OF THE PROPAGATION

Let us explicitate the second-order moment in channel c of \mathbf{t}^1 :

$$\mu_{2,c}(\mathbf{t}^1) = \mathbb{E}_{\mathbf{x},\mathbf{s},\alpha} \left[\hat{\mathbf{f}}_m(\mathbf{t}^1, \alpha)_c^2 \right] = \mathbb{E}_{\mathbf{x},\mathbf{s},\alpha} \left[\mathbf{f}_m(\mathbf{t}^1, \alpha)_c^2 \right] = \mathbb{E}_{\mathbf{x},\mathbf{s},\alpha} \left[(\mathbf{W}_{c,:}^1 \mathbf{r}_f(\mathbf{s}, \alpha))^2 \right] \quad (68)$$

$$= \sum_{i,j} \mathbf{W}_{c,i}^1 \mathbf{W}_{c,j}^1 \mathbb{E}_{\mathbf{s},\alpha} [\mathbf{r}_f(\mathbf{s}, \alpha)_i \mathbf{r}_f(\mathbf{s}, \alpha)_j] = \sum_i (\mathbf{W}_{c,i}^1)^2 = \|\mathbf{W}_{c,:}^1\|_2^2. \quad (69)$$

where Eq. (68) follows from \mathbf{t}^1 being centered and Eq. (69) follows from the white noise property $\mathbb{E}_{\mathbf{s}}[\mathbf{s}_i \mathbf{s}_j] = \delta_{ij}$, which implies for any α that $\mathbb{E}_{\mathbf{s}}[\mathbf{r}_f(\mathbf{s}, \alpha)_i \mathbf{r}_f(\mathbf{s}, \alpha)_j] = \delta_{ij}$ under periodic boundary conditions.

Now we turn to the second-order moment in channel c of \mathbf{y}^1 . Denoting $(\mathbf{e}_1, \dots, \mathbf{e}_{R_1})$ and $(\lambda_1, \dots, \lambda_{R_1})$ respectively the orthogonal eigenvectors and eigenvalues of $\mathbf{C}[\mathbf{r}_f(\mathbf{x}, \alpha)]$, and $\tilde{\mathbf{W}}^1 = \mathbf{W}^1(\mathbf{e}_1, \dots, \mathbf{e}_{R_1})$, we get

$$\begin{aligned} \mu_{2,c}(\mathbf{y}^1) &= \mathbb{E}_{\mathbf{x},\alpha} \left[\hat{\mathbf{f}}_m(\mathbf{y}^1, \alpha)_c^2 \right] = \mathbb{E}_{\mathbf{x},\alpha} \left[(\mathbf{W}_{c,:}^1 \hat{\mathbf{r}}_f(\mathbf{x}, \alpha))^2 \right] = \sum_i (\tilde{\mathbf{W}}_{c,i}^1)^2 \lambda_i \\ &= \|\mathbf{W}_{c,:}^1\|_2^2 \sum_i (\tilde{\mathbf{W}}_{c,i}^1)^2 \lambda_i = \mu_{2,c}(\mathbf{t}^1) \sum_i (\tilde{\mathbf{W}}_{c,i}^1)^2 \lambda_i, \end{aligned}$$

where we defined $\tilde{\mathbf{W}}^1$ such that $\forall c: \tilde{\mathbf{W}}_{c,:}^1 = \mathbf{W}_{c,:}^1 / \|\mathbf{W}_{c,:}^1\|$ and we used Eq. (69). Under standard initialization, the distribution of \mathbf{W}^1 is spherically symmetric, implying that for all c the distribution of $\tilde{\mathbf{W}}_{c,:}^1$ is uniform on the sphere of \mathbb{R}^{R_1} . In turn, this implies

$$\begin{aligned} \forall i: \mathbb{E}_{\theta^1} \left[(\tilde{\mathbf{W}}_{c,i}^1)^2 \right] &= \frac{1}{R_1}, \\ \forall c: \mathbb{E}_{\theta^1} \left[\sum_i (\tilde{\mathbf{W}}_{c,i}^1)^2 \lambda_i \right] &= \frac{1}{R_1} \sum_i \lambda_i, \quad \mathbb{E}_{c,\theta^1} \left[\sum_i (\tilde{\mathbf{W}}_{c,i}^1)^2 \lambda_i \right] = \frac{1}{R_1} \sum_i \lambda_i. \end{aligned} \quad (70)$$

Finally we can write $\exp(\bar{m}_{BN}[\zeta^1])$ as

$$\begin{aligned} \exp(\bar{m}_{BN}[\zeta^1]) &= \left(\frac{\mu_2(\mathbf{s}^0)}{\mu_2(\mathbf{x}^0)} \right)^{-1/2} \mathbb{E}_{c,\theta^1} \left[\frac{\mu_{2,c}(\mathbf{t}^1)}{\mu_{2,c}(\mathbf{y}^1)} \right]^{1/2} \\ &= \left(\frac{1}{R_1} \sum_i \lambda_i \right)^{1/2} \mathbb{E}_{c,\theta^1} \left[\frac{1}{\sum_i (\tilde{\mathbf{W}}_{c,i}^1)^2 \lambda_i} \right]^{1/2}, \end{aligned} \quad (71)$$

$$\geq \left(\frac{1}{R_1} \sum_i \lambda_i \right)^{1/2} \left(\mathbb{E}_{c,\theta^1} \left[\sum_i (\tilde{\mathbf{W}}_{c,i}^1)^2 \lambda_i \right]^{-1} \right)^{1/2} = 1. \quad (72)$$

where Eq. (71) is obtained using $\mu_2(\mathbf{s}^0) = 1$ and $\mu_2(\mathbf{x}^0) = \frac{1}{R_1} \text{Tr } \mathbf{C}[\mathbf{r}_f(\mathbf{x}, \boldsymbol{\alpha})] = \frac{1}{R_1} \sum_i \lambda_i$ by Corollary 3 while Eq. (72) is obtained using the convexity of $x \rightarrow 1/x$ and Eq. (70).

H NORMALIZED SENSITIVITY INCREMENTS OF BATCH-NORMALIZED RESNETS

H.1 ADAPTATION OF THE PREVIOUS SETUP TO RESNETS

Before proceeding to the analysis, slight adaptations and forewords are necessary. Let us denote $\Theta^{l,h} = (\boldsymbol{\omega}^{1,1}, \boldsymbol{\beta}^{1,1}, \dots, \boldsymbol{\omega}^{1,H}, \boldsymbol{\beta}^{1,H}, \dots, \boldsymbol{\omega}^{l,1}, \boldsymbol{\beta}^{l,1}, \dots, \boldsymbol{\omega}^{l,h}, \boldsymbol{\beta}^{l,h})$ for the full set of parameters up to layer h in residual unit l and $\theta^{l,h} = \Theta^{l,h} | \Theta^{l,h-1}$ for the conditional set of parameters of layer h in residual unit l . We further denote $\Theta^l = \Theta^{l,H}$ and $\theta^l = \Theta^{l,H} | \Theta^{l-1,H}$ respectively the full and conditional sets of parameters at the granularity of the residual unit.

We now clarify to what extent Theorem 3 on batch-normalized feedforward nets still apply. First let us rewrite the propagation at layer $1 \leq h \leq H$ inside residual unit l with the pre-activation perspective:

$$\mathbf{z}^{l,h} = BN(\mathbf{y}^{l,h-1}), \quad \mathbf{x}^{l,h} = \phi(\mathbf{z}^{l,h}), \quad \mathbf{y}^{l,h} = \boldsymbol{\omega}^{l,h} * \mathbf{x}^{l,h} + \boldsymbol{\beta}^{l,h}, \quad (73)$$

$$\mathbf{u}^{l,h} = BN'(\mathbf{y}^{l,h-1}) \odot \mathbf{t}^{l,h-1}, \quad \mathbf{s}^{l,h} = \phi'(\mathbf{z}^{l,h}) \odot \mathbf{u}^{l,h}, \quad \mathbf{t}^{l,h} = \boldsymbol{\omega}^{l,h} * \mathbf{s}^{l,h}, \quad (74)$$

In the pre-activation perspective, each layer starts with $(\mathbf{y}^{l,h-1}, \mathbf{t}^{l,h-1})$ after the convolution and ends at $(\mathbf{y}^{l,h}, \mathbf{t}^{l,h})$ again after the convolution. The concrete effect is that in the first layer $h = 1$ of each residual unit l , batch normalization and ϕ are completely deterministic conditionally on Θ^{l-1} . This occurs again for $h \geq 2$ since batch normalization and ϕ are random conditionally on Θ^{l-1} but completely deterministic conditionally on $\Theta^{l,h-1}$. At even larger granularity, due to the aggregation $(\mathbf{y}^l, \mathbf{t}^l) = \sum_{k=0}^l (\mathbf{y}^{k,H}, \mathbf{t}^{k,H})$, the input signal $(\mathbf{y}^{l-1}, \mathbf{t}^{l-1})$ of each residual unit becomes more and more correlated between successive l and less and less dependent on the parameters θ^{l-k} of individual previous units.

Since batch normalization and the nonlinearity ϕ are the drivers of the evolution of ζ^l , this shift can be thought as attributing the parameters and thus the stochasticity of layer h to the layer $h - 1$. A possible strategy is thus to consider the evolution from $(\mathbf{x}^{l,h-1}, \mathbf{s}^{l,h-1})$ to $(\mathbf{x}^{l,h}, \mathbf{s}^{l,h})$ for layers $2 \leq h \leq H$. This strategy however does not work for the first layer $h = 1$, since the input signal $(\mathbf{y}^{l-1}, \mathbf{t}^{l-1})$ depends on the whole sequence of parameters Θ^{l-1} and not only on $\theta^{l-1,H}$. Another strategy consists in treating all moment-related quantities as deterministic instead of random. Symmetric propagation does not occur strictly in this case, but it still occurs in a *mean-field sense* when averaging over channel. So we expect Theorem 3 to remain valid.

H.2 LEMMA ON DOT-PRODUCT

Lemma 13. For any random tensor \mathbf{u} of $\mathbb{R}^{n \times \dots \times n \times N}$:

$$\begin{aligned} \mathbb{E}_{\theta^l} [\mathbb{E}_{\mathbf{y}, \boldsymbol{\alpha}, \mathbf{c}} [\hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{u}, \boldsymbol{\alpha})_{\mathbf{c}} \hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{y}^{l,H}, \boldsymbol{\alpha})_{\mathbf{c}}]] &= 0, \\ \mathbb{E}_{\theta^l} [\mathbb{E}_{\mathbf{y}, \boldsymbol{\alpha}, \mathbf{c}} [\hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{u}, \boldsymbol{\alpha})_{\mathbf{c}} \hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{y}^{l,H}, \boldsymbol{\alpha})_{\mathbf{c}}]^2]] &\leq \frac{1}{N r_{\text{eff}}(\mathbf{u})} \mu_2(\mathbf{u}) \mathbb{E}_{\theta^l} [\mu_2(\mathbf{y}^{l,H})], \\ \mathbb{E}_{\theta^l} [\mathbb{E}_{\mathbf{y}, \mathbf{t}, \boldsymbol{\alpha}, \mathbf{c}} [\hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{u}, \boldsymbol{\alpha})_{\mathbf{c}} \hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{t}^{l,H}, \boldsymbol{\alpha})_{\mathbf{c}}]] &= 0, \\ \mathbb{E}_{\theta^l} [\mathbb{E}_{\mathbf{y}, \mathbf{t}, \boldsymbol{\alpha}, \mathbf{c}} [\hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{u}, \boldsymbol{\alpha})_{\mathbf{c}} \hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{t}^{l,H}, \boldsymbol{\alpha})_{\mathbf{c}}]^2]] &\leq \frac{1}{N r_{\text{eff}}(\mathbf{u})} \mu_2(\mathbf{u}) \mathbb{E}_{\theta^l} [\mu_2(\mathbf{t}^{l,H})]. \end{aligned}$$

Proof. By spherical symmetry, moments of $\hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{y}^{l,H}, \boldsymbol{\alpha})_{\mathbf{c}}$ and $-\hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{y}^{l,H}, \boldsymbol{\alpha})_{\mathbf{c}} = \hat{\mathbf{f}}_{\mathbf{m}}(-\mathbf{y}^{l,H}, \boldsymbol{\alpha})_{\mathbf{c}}$ have the same distribution with respect to θ^l . It follows that

$$\begin{aligned} \mathbb{E}_{\theta^l} [\mathbb{E}_{\mathbf{y}, \boldsymbol{\alpha}, \mathbf{c}} [\hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{u}, \boldsymbol{\alpha})_{\mathbf{c}} \hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{y}^{l,H}, \boldsymbol{\alpha})_{\mathbf{c}}]] &= \mathbb{E}_{\theta^l} [\mathbb{E}_{\mathbf{y}, \boldsymbol{\alpha}, \mathbf{c}} [\hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{u}, \boldsymbol{\alpha})_{\mathbf{c}} (-\hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{y}^{l,H}, \boldsymbol{\alpha})_{\mathbf{c}})]], \\ \mathbb{E}_{\theta^l} [\mathbb{E}_{\mathbf{y}, \boldsymbol{\alpha}, \mathbf{c}} [\hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{u}, \boldsymbol{\alpha})_{\mathbf{c}} \hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{y}^{l,H}, \boldsymbol{\alpha})_{\mathbf{c}}]] &= 0. \end{aligned}$$

Next we note that

$$\begin{aligned}\mathbb{E}_{\mathbf{y},\alpha,c}[\hat{\mathbf{f}}_m(\mathbf{u},\alpha)_c\hat{\mathbf{f}}_m(\mathbf{y}^{l,H},\alpha)_c] &= \frac{1}{N}\sum_c\mathbb{E}_{\mathbf{y},\alpha}[\hat{\mathbf{f}}_m(\mathbf{u},\alpha)_c\hat{\mathbf{f}}_m(\mathbf{y}^{l,H},\alpha)_c], \\ &= \frac{1}{N}\mathbb{E}_{\mathbf{y},\alpha}[\langle\hat{\mathbf{f}}_m(\mathbf{u},\alpha)\hat{\mathbf{f}}_m(\mathbf{y}^{l,H},\alpha)\rangle],\end{aligned}\quad (75)$$

where $\langle\cdot\rangle$ denotes the standard dot product in \mathbb{R}^N . Let us denote $(\mathbf{e}_1,\dots,\mathbf{e}_N)$ and $(\lambda_1,\dots,\lambda_N)$ respectively the orthogonal eigenvectors and eigenvalues of $\mathbf{C}[\hat{\mathbf{f}}_m(\mathbf{u},\alpha)]$. We further denote u_i the unit-variance components of $\hat{\mathbf{f}}_m(\mathbf{u},\alpha)$ in the basis $(\mathbf{e}_1,\dots,\mathbf{e}_N)$, and y_i the components of $\hat{\mathbf{f}}_m(\mathbf{y}^{l,H},\alpha)$ in the basis $(\mathbf{e}_1,\dots,\mathbf{e}_N)$. This gives

$$\begin{aligned}\hat{\mathbf{f}}_m(\mathbf{u},\alpha) &= \sum_i\sqrt{\lambda_i}u_i\mathbf{e}_i, \quad \mathbb{E}_{\mathbf{y},\alpha}[u_i] = 0, \quad \mathbb{E}_{\mathbf{y},\alpha}[u_i^2] = 1, \\ \hat{\mathbf{f}}_m(\mathbf{y}^{l,H},\alpha) &= \sum_i y_i\mathbf{e}_i.\end{aligned}$$

Now we decompose each component y_i of $\mathbf{y}^{l,H}$ as:

$$\forall j : \alpha_{i,j} = \mathbb{E}_{\mathbf{y},\alpha}[y_i u_j], \quad y_i = \sum_j \alpha_{i,j} u_j + z_i,$$

From this definition, we get

$$\begin{aligned}\forall j : \mathbb{E}_{\mathbf{y},\alpha}[z_i u_j] &= 0, \quad \mathbb{E}_{\mathbf{y},\alpha}[y_i u_i] = \alpha_{i,i}, \quad \mathbb{E}_{\mathbf{y},\alpha}[y_i^2] = \sum_j \alpha_{i,j}^2 + \mathbb{E}_{\mathbf{y},\alpha}[z_i^2], \\ \mu_2(\mathbf{y}^{l,H}) &= \frac{1}{N}\mathbb{E}_{\mathbf{y},\alpha}[\langle\mathbf{y}^{l,H},\mathbf{y}^{l,H}\rangle] = \frac{1}{N}\sum_i\mathbb{E}_{\mathbf{y},\alpha}[y_i^2] = \frac{1}{N}\left(\sum_{i,j}\alpha_{i,j}^2 + \sum_i\mathbb{E}_{\mathbf{y},\alpha}[z_i^2]\right).\end{aligned}\quad (76)$$

The dot product can be computed in any orthogonal basis, so we use the basis $(\mathbf{e}_1,\dots,\mathbf{e}_N)$:

$$\mathbb{E}_{\mathbf{y},\alpha}[\langle\hat{\mathbf{f}}_m(\mathbf{u},\alpha)\hat{\mathbf{f}}_m(\mathbf{y}^{l,H},\alpha)\rangle] = \sum_i\sqrt{\lambda_i}\mathbb{E}_{\mathbf{y},\alpha}[y_i u_i] = \sum_i\sqrt{\lambda_i}\alpha_{i,i}.$$

Spherical symmetry implies that moments of $y_1\mathbf{e}_1 + \dots + y_i\mathbf{e}_i + \dots + y_N\mathbf{e}_N$ and $y_1\mathbf{e}_1 + \dots - y_i\mathbf{e}_i + \dots + y_N\mathbf{e}_N$ have the same distribution with respect to θ^l . It follows that

$$\begin{aligned}\forall j \neq i : \mathbb{E}_{\mathbf{y},\alpha}[y_i u_i]\mathbb{E}_{\mathbf{y},\alpha}[y_j u_j] &\stackrel{\theta^l}{\sim} \mathbb{E}_{\mathbf{y},\alpha}[-y_i u_i]\mathbb{E}_{\mathbf{y},\alpha}[y_j u_j], \\ \forall j \neq i : \alpha_{i,i}\alpha_{j,j} &\stackrel{\theta^l}{\sim} (-\alpha_{i,i})\alpha_{j,j}, \\ \forall j \neq i : \mathbb{E}_{\theta^l}[\alpha_{i,i}\alpha_{j,j}] &= 0.\end{aligned}$$

We deduce that

$$\mathbb{E}_{\theta^l}\left[\mathbb{E}_{\mathbf{y},\alpha}[\langle\hat{\mathbf{f}}_m(\mathbf{u},\alpha)\hat{\mathbf{f}}_m(\mathbf{y}^{l,H},\alpha)\rangle]^2\right] = \sum_i\lambda_i\mathbb{E}_{\theta^l}[\alpha_{i,i}^2].$$

Spherical symmetry also implies that the distribution of $\alpha_{i,j}$ with respect to θ^l is the same for all i . Denoting (β_j) such that $\forall i, j : \beta_j = \mathbb{E}_{\theta^l}[\alpha_{i,j}^2]$, we get combined with Eq. (76):

$$\begin{aligned}\mathbb{E}_{\theta^l}[\mu_2(\mathbf{y}^{l,H})] &\geq \frac{1}{N}\left(\sum_{i,j}\beta_j\right) = \sum_i\beta_i, \\ \mathbb{E}_{\theta^l}\left[\mathbb{E}_{\mathbf{y},\alpha}[\langle\hat{\mathbf{f}}_m(\mathbf{u},\alpha)\hat{\mathbf{f}}_m(\mathbf{y}^{l,H},\alpha)\rangle]^2\right] &= \sum_i\lambda_i\beta_i \leq \lambda_{\max}\left(\sum_i\beta_i\right) \leq \lambda_{\max}\mathbb{E}_{\theta^l}[\mu_2(\mathbf{y}^{l,H})].\end{aligned}$$

Finally combining with Eq. (75):

$$\begin{aligned}\mathbb{E}_{\theta^l} \left[\mathbb{E}_{\mathbf{y}, \alpha, \mathbf{c}} \left[\hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{u}, \alpha)_{\mathbf{c}} \hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{y}^{l,H}, \alpha)_{\mathbf{c}} \right]^2 \right] &= \frac{1}{N^2} \mathbb{E}_{\theta^l} \left[\mathbb{E}_{\mathbf{y}, \alpha} \left[\langle \hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{u}, \alpha)_{\mathbf{c}} \hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{y}^{l,H}, \alpha)_{\mathbf{c}} \rangle \right]^2 \right] \\ &\leq \frac{1}{N^2} \lambda_{\max} \mathbb{E}_{\theta^l} [\mu_2(\mathbf{y}^{l,H})] \\ &\leq \frac{1}{Nr_{\text{eff}}(\mathbf{u})} \mu_2(\mathbf{u}) \mathbb{E}_{\theta^l} [\mu_2(\mathbf{y}^{l,H})],\end{aligned}$$

where we used $\lambda_{\max} r_{\text{eff}}(\mathbf{u}) = \sum_i \lambda_i = N \mu_2(\mathbf{u})$. The same analysis can be applied to $\hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{u}, \alpha)$ and $\hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{t}^{l,H}, \alpha)$. \square

Corollary 14. Let us denote the dot products for $k, l \geq 0$ as

$$Y_{k,l} = \mathbb{E}_{\mathbf{y}, \alpha, \mathbf{c}} [\hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{y}^{k,H}, \alpha)_{\mathbf{c}} \hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{y}^{l,H}, \alpha)_{\mathbf{c}}], \quad T_{k,l} = \mathbb{E}_{\mathbf{y}, \alpha, \mathbf{c}} [\hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{t}^{k,H}, \alpha)_{\mathbf{c}} \hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{t}^{l,H}, \alpha)_{\mathbf{c}}],$$

which combined with Eq. (18) implies

$$\begin{aligned}\sum_{k=0}^{l-1} Y_{k,l} &= \mathbb{E}_{\mathbf{y}, \alpha, \mathbf{c}} [\hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{y}^{l-1}, \alpha)_{\mathbf{c}} \hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{y}^{l,H}, \alpha)_{\mathbf{c}}], \\ \sum_{k=0}^{l-1} T_{k,l} &= \mathbb{E}_{\mathbf{y}, \alpha, \mathbf{c}} [\hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{t}^{l-1}, \alpha)_{\mathbf{c}} \hat{\mathbf{f}}_{\mathbf{m}}(\mathbf{t}^{l,H}, \alpha)_{\mathbf{c}}].\end{aligned}$$

Then by spherical symmetry $\forall k, l$, $Y_{k,l}$ and $T_{k,l}$ are centered, and $\forall \{k, l\} \neq \{k', l'\}$: $\mathbb{E}_{\Theta^{\max(k,l,k',l')}} [Y_{k,l} Y_{k',l'}] = 0$, $\mathbb{E}_{\Theta^{\max(k,l,k',l')}} [T_{k,l} T_{k',l'}] = 0$. Furthermore:

$$\forall k < l: \mathbb{E}_{\Theta^l} [Y_{k,l}^2] \leq \mathbb{E}_{\Theta^{l-1}} \left[\frac{1}{Nr_{\text{eff}}(\mathbf{y}^{k,H})} \mu_2(\mathbf{y}^{k,H}) \mathbb{E}_{\theta^l} [\mu_2(\mathbf{y}^{l,H})] \right], \quad (77)$$

$$\forall k < l: \mathbb{E}_{\Theta^l} [T_{k,l}^2] \leq \mathbb{E}_{\Theta^{l-1}} \left[\frac{1}{Nr_{\text{eff}}(\mathbf{t}^{k,H})} \mu_2(\mathbf{t}^{k,H}) \mathbb{E}_{\theta^l} [\mu_2(\mathbf{t}^{l,H})] \right], \quad (78)$$

$$\forall l: \mathbb{E}_{\Theta^l} \left[\left(\sum_{k=0}^{l-1} Y_{k,l} \right)^2 \right] \leq \mathbb{E}_{\Theta^{l-1}} \left[\frac{1}{Nr_{\text{eff}}(\mathbf{y}^{l-1})} \mu_2(\mathbf{y}^{l-1}) \mathbb{E}_{\theta^l} [\mu_2(\mathbf{y}^{l,H})] \right], \quad (79)$$

$$\forall l: \mathbb{E}_{\Theta^l} \left[\left(\sum_{k=0}^{l-1} T_{k,l} \right)^2 \right] \leq \mathbb{E}_{\Theta^{l-1}} \left[\frac{1}{Nr_{\text{eff}}(\mathbf{t}^{l-1})} \mu_2(\mathbf{t}^{l-1}) \mathbb{E}_{\theta^l} [\mu_2(\mathbf{t}^{l,H})] \right]. \quad (80)$$

H.3 PROOF OF THEOREM 4

Theorem 4. Normalized Sensitivity increments of batch-normalized resnets. Suppose that for all depth l we can bound the effective ranks $r_{\min} \lesssim r_{\text{eff}}(\mathbf{y}^l), r_{\text{eff}}(\mathbf{y}^{l,H}), r_{\text{eff}}(\mathbf{t}^l), r_{\text{eff}}(\mathbf{t}^{l,H})$, the second-order central moment $\mu_{2,\min} \lesssim \mu_2(\mathbf{y}^{l,H}) \lesssim \mu_{2,\max}$ and the feedforward increments inside residual units $\delta_{\min} \lesssim \delta \zeta^{l,h} \lesssim \delta_{\max}$. Denote $\rho_{\min} = ((\delta_{\min})^{2H} \mu_{2,\min} - \mu_{2,\max}) / \mu_{2,\max}$ and $\rho_{\max} = ((\delta_{\max})^{2H} \mu_{2,\max} - \mu_{2,\min}) / \mu_{2,\min}$, and further consider τ_{\min}, τ_{\max} such that $\tau_{\min} < \rho_{\min} / 2$ and $\tau_{\max} > \rho_{\max} / 2$. Then:

$$\forall l \ll Nr_{\min}: \left(1 + \frac{\rho_{\min}}{l+1} \right)^{1/2} \lesssim \delta \zeta^l \lesssim \left(1 + \frac{\rho_{\max}}{l+1} \right)^{1/2}, \quad (81)$$

$$\forall l \gg 1: \frac{1}{2} \rho_{\min} \log l \lesssim \log \zeta^l \lesssim \frac{1}{2} \rho_{\max} \log l, \quad (82)$$

$$\forall l \gg 1: l^{\tau_{\min}} \lesssim \zeta^l \lesssim l^{\tau_{\max}}. \quad (83)$$

Proof. First we introduce the additional constants $\gamma_{\min} = (\delta_{\min})^{2H}$ and $\gamma_{\max} = (\delta_{\max})^{2H}$, so that we can write $\rho_{\min} = (\gamma_{\min}\mu_{2,\min} - \mu_{2,\max}) / \mu_{2,\max}$ and $\rho_{2,\max} = (\gamma_{\max}\mu_{2,\max} - \mu_{2,\min}) / \mu_{2,\min}$.

We also remind that the symbols \lesssim in Theorem 4 denote inequalities up to small non-dominating terms. We write $a \lesssim b$ when $a(1 + \delta_a) \leq b(1 + \delta_b)$ with $|\delta_a| \ll 1$, $|\delta_b| \ll 1$ with high probability. We write $a \simeq b$ when $a(1 + \delta_a) = b(1 + \delta_b)$ with $|\delta_a| \ll 1$, $|\delta_b| \ll 1$ with high probability. Denoting \wedge for the logical *and*, the following rules are easily verified:

$$\begin{aligned} (a \lesssim b) \wedge (a \gtrsim b) &\iff (a \simeq b), \\ (a \lesssim b) &\implies (-a \gtrsim -b), \\ (a \lesssim b) &\implies (1/a \gtrsim 1/b), \\ (a \lesssim b) \wedge (c \lesssim d) &\implies (ac \lesssim bd), \\ (a \lesssim b) \wedge (b \lesssim c) &\implies a \lesssim c. \end{aligned}$$

Finally $(a \lesssim b) \wedge (c \lesssim d) \implies (a + c \lesssim b + d)$ under the condition that $\{|a + c| \ll |a| + |c|\}$ and $\{|b + d| \ll |b| + |d|\}$ are very small probability events. We keep these rules in mind in the course of this proof.

Proof of Eq. (81). Adopting the same notations as Corollary 14 and using $\mathbf{y}^l = \sum_{k=0}^l \mathbf{y}^{k,H}$ by Eq. (18), we get

$$\begin{aligned} \mu_2(\mathbf{y}^l) &= \mathbb{E}_{\mathbf{y}, \alpha, c} \left[\sum_{k, k'} \hat{\mathbf{f}}_m(\mathbf{y}^{k,H}, \alpha)_c \hat{\mathbf{f}}_m(\mathbf{y}^{k',H}, \alpha)_c \right] = \sum_{k, k'} Y_{k, k'}, \\ \mu_2(\mathbf{t}^l) &= \mathbb{E}_{\mathbf{y}, \alpha, c} \left[\sum_{k, k'} \hat{\mathbf{f}}_m(\mathbf{t}^{k,H}, \alpha)_c \hat{\mathbf{f}}_m(\mathbf{t}^{k',H}, \alpha)_c \right] = \sum_{k, k'} T_{k, k'}. \end{aligned}$$

Now using the hypotheses $\mu_{2,\min} \lesssim \mu_2(\mathbf{y}^{l,H}) \lesssim \mu_{2,\max}$ and $r_{\min} \lesssim r_{\text{eff}}(\mathbf{y}^{l,H})$, combined with Eq. (77) from Corollary 14,

$$\begin{aligned} (l+1)\mu_{2,\min} + \sum_{k \neq k'} Y_{k, k'} &\lesssim \mu_2(\mathbf{y}^l) \lesssim (l+1)\mu_{2,\max} + \sum_{k \neq k'} Y_{k, k'}, \\ \mathbb{E}_{\Theta^l} \left[\left(\sum_{k \neq k'} Y_{k, k'} \right)^2 \right] &\lesssim l(l+1) \frac{1}{Nr_{\min}} \mu_{2,\max}^2, \\ \mathbb{E}_{\Theta^l} \left[\left| \sum_{k \neq k'} Y_{k, k'} \right| \right] &\lesssim (l+1) \frac{1}{\sqrt{Nr_{\min}}} \mu_{2,\max}, \end{aligned} \quad (84)$$

where Eq. (84) is obtained using Cauchy-Schwarz inequality. It follows that for large width $N \gg 1$, with high probability $\left| \sum_{k \neq k'} Y_{k, k'} \right| \ll (l+1)\mu_{2,\min}$ and $\left| \sum_{k \neq k'} Y_{k, k'} \right| \ll (l+1)\mu_{2,\max}$, which then gives

$$(l+1)\mu_{2,\min} \lesssim \mu_2(\mathbf{y}^l) \lesssim (l+1)\mu_{2,\max}. \quad (85)$$

We can write $(\zeta^l)^2$ as

$$\begin{aligned} (\zeta^l)^2 &= \frac{\mu_2(\mathbf{y}^0)\mu_2(\mathbf{t}^l)}{\mu_2(\mathbf{y}^l)} = \mu_2(\mathbf{y}^0) \frac{\mu_2(\mathbf{t}^{l-1}) + T_{l,l} + 2 \sum_{k=0}^{l-1} T_{k,l}}{\mu_2(\mathbf{y}^{l-1}) + Y_{l,l} + 2 \sum_{k=0}^{l-1} Y_{k,l}}, \\ (\zeta^l)^2 &= (\zeta^{l-1})^2 \frac{\mu_2(\mathbf{y}^{l-1}) + \frac{\mu_2(\mathbf{y}^0)}{(\zeta^{l-1})^2} T_{l,l} + 2 \frac{\mu_2(\mathbf{y}^0)}{(\zeta^{l-1})^2} \sum_{k=0}^{l-1} T_{k,l}}{\mu_2(\mathbf{y}^{l-1}) + Y_{l,l} + 2 \sum_{k=0}^{l-1} Y_{k,l}}. \end{aligned} \quad (86)$$

Let us denote $\forall k \leq l$: $\tilde{T}_{k,l} = \frac{\mu_2(\mathbf{y}^0)}{(\zeta^{l-1})^2} T_{k,l}$ and $Y_l = \sum_{k=0}^{l-1} Y_{k,l}$ and $\tilde{T}_l = \sum_{k=0}^{l-1} \tilde{T}_{k,l}$. Eq. (86) then gives

$$(\delta\zeta^l)^2 = \frac{(\zeta^l)^2}{(\zeta^{l-1})^2} = \frac{\mu_2(\mathbf{y}^{l-1}) + \tilde{T}_{l,l} + 2\tilde{T}_l}{\mu_2(\mathbf{y}^{l-1}) + Y_{l,l} + 2Y_l}. \quad (87)$$

Furthermore we can bound $\tilde{T}_{l,l}$ as

$$\begin{aligned} \tilde{T}_{l,l} &= \frac{\mu_2(\mathbf{y}^0)}{(\zeta^{l-1})^2} \mu_2(\mathbf{t}^{l,H}) = \frac{\mu_2(\mathbf{y}^0)}{(\zeta^{l-1})^2} (\zeta^{l-1})^2 \prod_h (\delta\zeta^{l,h})^2 \frac{\mu_2(\mathbf{y}^{l,H})}{\mu_2(\mathbf{y}^0)}, \\ \gamma_{\min}\mu_{2,\min} &\lesssim \tilde{T}_{l,l} \lesssim \gamma_{\max}\mu_{2,\max}. \end{aligned} \quad (88)$$

Combining Eq. (79) and Eq. (80) from Corollary 14 with Eq. (85), we get for the variance of the terms \tilde{T}_l and Y_l :

$$\mathbb{E}_{\Theta^l} [Y_l^2] \lesssim \frac{1}{Nr_{\min}} l \mu_{2,\max}^2, \quad (89)$$

$$\begin{aligned} \mathbb{E}_{\Theta^l} \left[\left(\frac{\mu_2(\mathbf{y}^0)}{(\zeta^{l-1})^2} \sum_{k < l} T_{k,l} \right)^2 \right] &\lesssim \frac{1}{Nr_{\min}} \mathbb{E}_{\Theta^{l-1}} \left[\frac{\mu_2(\mathbf{y}^0) \mu_2(\mathbf{t}^{l-1})}{(\zeta^{l-1})^2} \mathbb{E}_{\Theta^l} \left[\frac{\mu_2(\mathbf{y}^0) \mu_2(\mathbf{t}^{l,H})}{(\zeta^{l-1})^2} \right] \right], \\ \mathbb{E}_{\Theta^l} [\tilde{T}_l^2] &\lesssim \frac{1}{Nr_{\min}} \mathbb{E}_{\Theta^{l-1}} [\mu_2(\mathbf{y}^{l-1})] \mathbb{E}_{\Theta^l} [\tilde{T}_{l,l}] \lesssim \gamma_{\max} \frac{1}{Nr_{\min}} l \mu_{2,\max}^2. \end{aligned} \quad (90)$$

It follows that $|Y_l| \ll 1$ and $|\tilde{T}_l| \ll 1$ with high probability when $l \ll Nr_{\min}$. Combined with Eq. (87) and Eq. (88),

$$\begin{aligned} \frac{\mu_2(\mathbf{y}^{l-1}) + \gamma_{\min} \mu_{2,\min}}{\mu_2(\mathbf{y}^{l-1}) + \mu_{2,\max}} &\lesssim (\delta \zeta^l)^2 \lesssim \frac{\mu_2(\mathbf{y}^{l-1}) + \gamma_{\max} \mu_{2,\max}}{\mu_2(\mathbf{y}^{l-1}) + \mu_{2,\min}}, \\ 1 + \frac{\gamma_{\min} \mu_{2,\min} - \mu_{2,\max}}{(l+1) \mu_{2,\max}} &\lesssim (\delta \zeta^l)^2 \lesssim 1 + \frac{\gamma_{\max} \mu_{2,\max} - \mu_{2,\min}}{(l+1) \mu_{2,\min}}. \end{aligned}$$

Finally for $l \ll Nr_{\min}$, we find that $\delta \zeta^l$ is bounded as

$$\left(1 + \frac{\rho_{\min}}{l+1}\right)^{1/2} \lesssim \delta \zeta^l \lesssim \left(1 + \frac{\rho_{\max}}{l+1}\right)^{1/2}.$$

□

Proof of Eq. (82). Using Eq. (89) and Eq. (90), we apply Cauchy-Schwarz inequality on Y_l and \tilde{T}_l :

$$\begin{aligned} \mathbb{E}_{\Theta^l} [|Y_l|] &\lesssim \frac{1}{\sqrt{Nr_{\min}}} \sqrt{l} \mu_{2,\max}, \\ \mathbb{E}_{\Theta^l} [|\tilde{T}_l|] &\lesssim \sqrt{\gamma_{\max}} \frac{1}{\sqrt{Nr_{\min}}} \sqrt{l} \mu_{2,\max}. \end{aligned}$$

Combined with Eq. (85), we deduce that for large width $|Y_l| \ll \mu_2(\mathbf{y}^l)$ and $|\tilde{T}_l| \ll \mu_2(\mathbf{y}^l)$ always hold with high probability. Due to $\mu_{2,\min} \lesssim Y_{l,l} \lesssim \mu_{2,\max}$ and to Eq. (88), we also have with high probability that $Y_{l,l} \ll \mu_2(\mathbf{y}^l)$ and $\tilde{T}_l \ll \mu_2(\mathbf{y}^l)$ when $l \gg 1$. Now let us take the logarithm in Eq. (87):

$$2 \log \delta \zeta^l = \log \left(1 + \frac{1}{\mu_2(\mathbf{y}^{l-1})} \tilde{T}_{l,l} + \frac{2}{\mu_2(\mathbf{y}^{l-1})} \tilde{T}_l \right) - \log \left(1 + \frac{1}{\mu_2(\mathbf{y}^{l-1})} Y_{l,l} + \frac{2}{\mu_2(\mathbf{y}^{l-1})} Y_l \right).$$

When $l \gg 1$, all the terms added to 1 in the logarithm are $\ll 1$ with high probability. Therefore for $l > l_0 \gg 1$, we can write

$$\begin{aligned} 2 \log \delta \zeta^l &\simeq \frac{1}{\mu_2(\mathbf{y}^{l-1})} \tilde{T}_{l,l} + \frac{2}{\mu_2(\mathbf{y}^{l-1})} \tilde{T}_l - \frac{1}{\mu_2(\mathbf{y}^{l-1})} Y_{l,l} - \frac{2}{\mu_2(\mathbf{y}^{l-1})} Y_l, \\ 2 \sum_{k=l_0+1}^l \log \delta \zeta^k &\simeq \sum_{k=l_0+1}^l \frac{1}{\mu_2(\mathbf{y}^{k-1})} (\tilde{T}_{k,k} - Y_{k,k}) + \sum_{k=l_0+1}^l \frac{2}{\mu_2(\mathbf{y}^{k-1})} (\tilde{T}_k - Y_k). \end{aligned}$$

Let us bound the first term:

$$\begin{aligned} \sum_{k=l_0+1}^l \frac{\gamma_{\min} \mu_{2,\min} - \mu_{2,\max}}{k \mu_{2,\max}} &\lesssim \sum_{k=l_0+1}^l \frac{1}{\mu_2(\mathbf{y}^{k-1})} (\tilde{T}_{k,k} - Y_{k,k}) \lesssim \sum_{k=l_0+1}^l \frac{\gamma_{\max} \mu_{2,\max} - \mu_{2,\min}}{k \mu_{2,\min}}, \\ \rho_{\min} \int_{l_0}^l \frac{1}{x} dx &\lesssim \sum_{k=l_0+1}^l \frac{1}{\mu_2(\mathbf{y}^{k-1})} (\tilde{T}_{k,k} - Y_{k,k}) \lesssim \rho_{\max} \int_{l_0}^l \frac{1}{x} dx, \\ \rho_{\min} \log \left(\frac{l}{l_0} \right) &\lesssim \sum_{k=l_0+1}^l \frac{1}{\mu_2(\mathbf{y}^{k-1})} (\tilde{T}_{k,k} - Y_{k,k}) \lesssim \rho_{\max} \log \left(\frac{l}{l_0} \right). \end{aligned} \quad (91)$$

Now we consider the second term. By spherical symmetry, $Y_k / \mu_2(\mathbf{y}^{k-1})$ and $Y_{k'} / \mu_2(\mathbf{y}^{k'-1})$ are non-correlated for $k \neq k'$. So we get combined with Eq. (89) that for $l > l_0 \gg 1$,

$$\begin{aligned} \mathbb{E}_{\Theta^l} \left[\left(\sum_{k=l_0+1}^l \frac{1}{\mu_2(\mathbf{y}^{k-1})} Y_k \right)^2 \right] &= \sum_{k=l_0+1}^l \mathbb{E}_{\Theta^k} \left[\left(\frac{1}{\mu_2(\mathbf{y}^{k-1})} Y_k \right)^2 \right], \\ &\lesssim \frac{\mu_{2,\max}^2}{Nr_{\min}\mu_{2,\min}^2} \sum_{k=l_0+1}^l \frac{k}{k^2} \lesssim \frac{\mu_{2,\max}^2}{Nr_{\min}\mu_{2,\min}^2} \log \left(\frac{l}{l_0} \right). \end{aligned}$$

A similar calculation for \tilde{T}_k gives

$$\mathbb{E}_{\Theta^l} \left[\left(\sum_{k=l_0+1}^l \frac{1}{\mu_2(\mathbf{y}^{k-1})} \tilde{T}_k \right)^2 \right] \lesssim \gamma_{\max} \frac{\mu_{2,\max}^2}{Nr_{\min}\mu_{2,\min}^2} \log \left(\frac{l}{l_0} \right).$$

For large width and $l \gg 1$, these terms are very small with high probability compared to the term of Eq. (91). It follows that for $l > l_0 \gg 1$,

$$\begin{aligned} \rho_{\min} \log \left(\frac{l}{l_0} \right) &\lesssim 2 \log \left(\frac{\zeta^l}{\zeta^{l_0}} \right) \lesssim \rho_{\max} \log \left(\frac{l}{l_0} \right), \\ \frac{1}{2} \rho_{\min} \log l &\lesssim \log \zeta^l \lesssim \frac{1}{2} \rho_{\max} \log l. \end{aligned}$$

□

Proof of Eq. (83). As a consequence Eq. (82), for $l \gg 1$ there exist δ, δ' with $|\delta| \ll 1, |\delta'| \ll 1$ with high probability and

$$\begin{aligned} (1 + \delta) \frac{1}{2} \rho_{\min} \log l &\leq \log \zeta^l \leq (1 + \delta') \frac{1}{2} \rho_{\max} \log l, \\ \exp \left(\frac{1}{2} (1 + \delta) \rho_{\min} \log l \right) &\leq \zeta^l \leq \exp \left(\frac{1}{2} (1 + \delta') \rho_{\max} \log l \right). \end{aligned}$$

Now consider τ_{\min}, τ_{\max} such that $\tau_{\min} < \frac{1}{2} \rho_{\min}$ and $\tau_{\max} > \frac{1}{2} \rho_{\max}$. Then

$$\begin{aligned} \frac{1}{2} (1 + \delta) \rho_{\min} \log l &= \left(\frac{1}{2} \rho_{\min} + \frac{1}{2} \rho_{\min} \delta - \tau_{\min} \right) \log l + \tau_{\min} \log l, \\ \frac{1}{2} (1 + \delta') \rho_{\max} \log l &= \left(\frac{1}{2} \rho_{\max} + \frac{1}{2} \rho_{\max} \delta' - \tau_{\max} \right) \log l + \tau_{\max} \log l. \end{aligned}$$

The terms $\frac{1}{2} \rho_{\min} + \frac{1}{2} \rho_{\min} \delta - \tau_{\min}$ and $\frac{1}{2} \rho_{\max} + \frac{1}{2} \rho_{\max} \delta' - \tau_{\max}$ are respectively positive with high probability and negative with high probability. Therefore with high probability, $\exp(\tau_{\min} \log l) \leq \zeta^l \leq \exp(\tau_{\max} \log l)$ and thus for $l \gg 1$:

$$\begin{aligned} \exp(\tau_{\min} \log l) &\lesssim \zeta^l \lesssim \exp(\tau_{\max} \log l), \\ l^{\tau_{\min}} &\lesssim \zeta^l \lesssim l^{\tau_{\max}}. \end{aligned}$$

□