



**HAL**  
open science

# Imposing Gaussian Pre-Activations in a Neural Network

Pierre Wolinski, Julyan Arbel

► **To cite this version:**

Pierre Wolinski, Julyan Arbel. Imposing Gaussian Pre-Activations in a Neural Network. JDS 2022 - 53es Journées de Statistique de la Société Française de Statistiques (SFdS), Jun 2022, Lyon, France. hal-03853790

**HAL Id: hal-03853790**

**<https://hal.science/hal-03853790>**

Submitted on 15 Nov 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Imposing Gaussian Pre-Activations in a Neural Network

Pierre Wolinski<sup>1</sup> and Julyan Arbel<sup>1</sup>

<sup>1</sup>Univ. Grenoble Alpes, Inria, CNRS, Grenoble INP, LJK, 38000 Grenoble, France, {pierre.wolinski|julyan.arbel}@inria.fr

May 31, 2022

## Abstract

The goal of the present work is to propose a way to modify both the initialization distribution of the weights of a neural network and its activation function, such that all pre-activations are Gaussian. We propose a family of pairs initialization/activation, where the activation functions span a continuum from bounded functions (such as Heaviside or tanh) to the identity function.

This work is motivated by the contradiction between existing works dealing with Gaussian pre-activations: on one side, the works in the line of the Neural Tangent Kernels and the Edge of Chaos are assuming it, while on the other side, theoretical and experimental results challenge this hypothesis.

The family of pairs initialization/activation we are proposing will help us to answer this hot question: is it desirable to have Gaussian pre-activations in a neural network?

## TODO:

- give more information on  $\phi_{2+}$

## 1 Introduction

Let us take a neural network at initialization with  $L$  layers. The operation made by the  $l$ -th layer is:

$$\forall l \in [1, L] : \quad X^{(l+1)} = \phi(Z^{(l+1)}) \quad (1)$$

$$\text{where:} \quad Z^{(l+1)} = \frac{1}{\sqrt{n_l}} W^{(l)} X^{(l)} + b^{(l)}, \quad (2)$$

where  $X^{(l+1)} \in \mathbb{R}^{n_{l+1}}$  is its *activation*,  $Z^{(l+1)} \in \mathbb{R}^{n_{l+1}}$  its *pre-activation*,  $\phi$  is the coordinate-wise *activation function*,  $W^{(l)} \in \mathbb{R}^{n_{l+1} \times n_l}$  is the *weight matrix* of the layer,  $b^{(l)} \in \mathbb{R}^{n_{l+1}}$  its *vector of biases*, and  $X^{(l)} \in \mathbb{R}^{n_l}$  its *input*.

The object of the present study is the distribution of the pre-activations  $Z^{(l)}$ , for a fixed input  $X^{(0)}$ , and weights  $W^{(l)}$  and biases  $b^{(l)}$  randomly initialized according to given distributions.

**Gaussian hypothesis for the pre-activations.** In the theoretical analysis of the properties of neural networks at initialization, the hypothesis of Gaussian pre-activations is common. In particular, this is a fundamental assumption when studying “Neural Tangent Kernels” (NTK) [3] or “Edge of Chaos” (EOC) [5]. In a nutshell, the NTK is an operator describing the optimization trajectory of an *infinitely wide* neural network (NN), which is believed to help to understand the optimization of ordinary NNs; the EOC is a criterion over the initialization distribution of the weights and biases of an *infinitely wide* NN, ensuring that information propagates and backpropagates the best across the layers of the NN. On one side, this *Gaussian hypothesis* can be justified in the case of “infinitely wide” NNs (i.e., when the widths  $n_l$  of the layers tend to infinity), by application of the Central Limit Theorem (see Eqn. (2) when  $n_l \rightarrow \infty$ ) [4]. On the other side, it is apparently necessary to get the results of [3, 5]. However, this Gaussian hypothesis remains debated for both theoretical and practical reasons.

First, from a strictly theoretical point of view, it has been shown that, for finite-width NNs (finite  $n_l$ ), the distribution of  $Z^{(l)}$  tends to have heavier and heavier tails as  $l$  increases, that is, as information flows from the input to the output [9]. Second, a series of experiments tends to show that pushing the distribution of the pre-activations towards a Gaussian (e.g., through a Bayesian prior) leads to worse performances than pushing it towards distributions with heavier tails (e.g., Laplace distribution) [?].

**Contributions.** The goal of the present work is testing the accuracy of neural networks with a *realistic* architecture, that is, of finite width ( $n_l$  may be small) and commonly used (“fully connected” or “convolutional”), *adjusted in such a way that all the pre-activations are Gaussian*. To this end, we propose a family of pairs  $(P_\theta, \phi_\theta)$ , where  $P_\theta$  is the distribution of the weights at initialization,  $\phi_\theta$  is the activation function, and  $\theta \in (2, \infty)$  is a parameter such that:

- for a given layer  $l$ :

$$\left. \begin{array}{l} Z_i^{(l)} \sim \mathcal{N}(0, 1) \text{ i.i.d.} \\ W_{ij}^{(l)} \sim P_\theta \text{ i.i.d.} \\ b_i^{(l)} = 0 \end{array} \right\} \Rightarrow Z_i^{(l+1)} := \frac{1}{\sqrt{n_l}} W^{(l)} \phi_\theta(Z^{(l)}) + b^{(l)} \sim \mathcal{N}(0, 1),$$

in other words, the pre-activations  $Z_i^{(l+1)}$  remain Gaussian for all  $l$ ;

- $P_\theta$  is the symmetric Weibull distribution  $\mathcal{W}(\theta, 1)$ , of CDF:

$$F_W(t) = \frac{1}{2} + \frac{1}{2} \operatorname{sgn}(t) \exp(-|t|^\theta); \quad (3)$$

- $\phi_\theta$  is defined in such a way that  $Z_i^{(l+1)} := \frac{1}{\sqrt{n_l}} W^{(l)} \phi_\theta(Z^{(l)}) + b^{(l)}$  is Gaussian;

- limiting cases:

$$\begin{aligned} \theta \rightarrow \infty &\Rightarrow P_\theta \xrightarrow{d} \mathcal{R} & \text{and} & \phi_\theta \xrightarrow{p.p.} \text{Id}, \\ \theta \rightarrow 2 &\Rightarrow P_\theta \xrightarrow{d} \mathcal{W}(2, 1) & \text{and} & \phi_\theta \xrightarrow{p.p.} \phi_{2+}, \end{aligned}$$

where  $\mathcal{R}$  is the Rademacher distribution (if  $\xi \sim \mathcal{R}$ , then  $\mathbb{P}(\xi = \pm 1) = 1/2$ ), Id is the identity function, and  $\phi_{2+}$  is an increasing function with:  $\lim_{\pm\infty} \phi_{2+} = \pm 1$  and  $\phi'_{2+}(0) = \sqrt{\pi}$ .

In the limiting case  $\theta \rightarrow \infty$ , we initialize the weights at  $\pm 1$ , which corresponds to binary “weight quantization”, used in the field of neural networks compression [6], and we use a linear activation function, commonly used in theoretical analyses of neural networks [1].

In the limiting case  $\theta \rightarrow 2$ , we recover weights whose distribution has a Gaussian tail, with an Heaviside-like activation function, used since the very beginning of the neural networks [7].

## 2 Obtaining the family $\{(P_\theta, \phi_\theta) : \theta \in (2, \infty)\}$

### 2.1 Decomposing the problem

In order to evaluate the distribution of the pre-activations outputted by a layer, we should be able to deal with a linear combination of products of random variables:

$$Z^{(l+1)} = \frac{1}{\sqrt{n_l}} W^{(l)} \phi(Z^{(l)}) + b^{(l)}. \quad (4)$$

For i.i.d. pre-activations  $Z_i^{(l)} \sim \mathcal{N}(0, 1)$ , we want the components of  $Z^{(l+1)}$  to be  $\mathcal{N}(0, 1)$ . Without loss of generality, nous consider only one component of  $Z^{(l+1)}$ , we omit layer indices, we replace  $Z^{(l)}$  by  $X$ , and the bias  $b^{(l)}$  by 0 (see discussion in Sec. 3):

$$Z = \frac{1}{\sqrt{n}} \sum_{i=1}^n W_i \phi(X_i).$$

Since we want to have  $Z \sim \mathcal{N}(0, 1)$ , and provided that  $Z$  is a linear combination of i.i.d. r.v.  $(W_i \phi(X_i))_i$ , we want to construct an initialization distribution for  $W_i$  and a function  $\phi$  in such a way that all  $(W_i \phi(X_i))_i$  follow a distribution  $\mathcal{N}(0, 1)$ . This is imposed by the following lemma:

**Lemma 1.** *Let  $(Z_i)_i$  be i.i.d. random variables. Let  $Z = \frac{1}{\sqrt{n}} \sum_{i=1}^n Z_i$ . If  $Z$  is  $\mathcal{N}(0, 1)$ , then the distribution of the  $Z_i$  is also  $\mathcal{N}(0, 1)$ .*

See proof in App. A.

**Main problem.** Let  $X \sim \mathcal{N}(0, 1)$ . We want to find a family  $\Theta$  of parameters  $\theta$  such that, for any  $\theta \in \Theta$ , there exists a probability distribution  $P_\theta$  and a function  $\phi_\theta$  such that:  $W \sim P_\theta \Rightarrow W\phi_\theta(X) \sim \mathcal{N}(0, 1)$ .

We decompose the problem into two parts, by introducing an intermediary random variable  $Y = \phi(X)$ :

- for a distribution  $P_\theta$ , deduce  $Q_\theta$  s.t.:  $W \sim P_\theta, Y \sim Q_\theta \Rightarrow WY =: G \sim \mathcal{N}(0, 1)$ ;
- for a distribution  $Q_\theta$ , find a function  $\phi_\theta$  s.t.:  $X \sim \mathcal{N}(0, 1) \Rightarrow Y = \phi_\theta(X) \sim Q_\theta$ .

## 2.2 Why choosing the family of Weibull distributions?

We are looking for a family of distributions  $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$  such that, for any  $\theta \in \Theta$ , there exists  $Q_\theta$  such that:

$$W \sim P_\theta, Y \sim Q_\theta \Rightarrow WY = G \sim \mathcal{N}(0, 1).$$

Therefore, the family  $\mathcal{P}$  is subject to several conditions. Two results tend to indicate that a subset of the family of Weibull distributions is a good choice:

1. according to the results of Section 2.2.1, the density of  $W$  at 0 should be 0, which is the case for all Weibull distributions;
2. according to the results of Section 2.2.2,  $W$  should be generalized Weibull-tail of parameter  $\theta \in (2, \infty)$  (see Definition 2).

In the process, we are able to gather information about the distribution of  $|Y|$ , namely its density at 0 and the leading power of the log of its survival function:

$$f_{|Y|}(0) = \sqrt{\frac{2}{\pi}} \left[ \int_0^\infty \frac{f_{|W|}(t)}{t} dt \right]^{-1},$$

$$\log S_{|Y|}(y) \propto -y^{1/(\frac{1}{2}-\frac{1}{\theta})}.$$

As a conclusion of this section, we consider that the distribution  $P_\theta$  of  $W$  lies in the following subset of the family of symmetric Weibull distributions (defined at Eqn. (3)):

$$\mathcal{P} := \{\mathcal{W}(\theta, 1) : \theta \in \Theta\},$$

$$\Theta := (2, \infty).$$

### 2.2.1 Behavior near 0

Since the product  $G = WY$  is  $\mathcal{N}(0, 1)$ , then we must have  $f_{|G|}(0) = \sqrt{2/\pi} \in (0, \infty)$ , which is impossible for several choices of distributions for  $W$ .

**Proposition 1** (density of a product of random variables at 0). *Let  $W, Y$  be two independent non-negative random variables and  $Z = WY$ . Let  $f_W, f_Y, f_Z$  be their respective density. Assuming that  $f_Y$  is continuous at 0 with  $f_Y(0) > 0$ , we have:*

$$\lim_{w \rightarrow 0} \int_w^\infty \frac{f_W(t)}{t} dt = \infty \quad \Rightarrow \quad \lim_{z \rightarrow 0} f_Z(z) = \infty, \quad (5)$$

moreover, if  $f_Y$  is bounded:

$$\int_0^\infty \frac{f_W(t)}{t} dt < \infty \quad \Rightarrow \quad f_Z(0) = f_Y(0) \int_0^\infty \frac{f_W(t)}{t} dt.$$

**Corollary 1.** *If  $f_Y$  and  $f_W$  are continuous at 0 with  $f_Y(0) > 0$  and  $f_Z(0) > 0$ , then:*

$$\lim_{z \rightarrow 0} f_Z(z) = \infty. \quad (6)$$

According to Corollary 1, it is impossible to obtain a Gaussian  $G$  by multiplying two random variables  $W$  and  $Y$  whose densities are both continuous and non-zero at 0. So, if we want to manipulate continuous densities, it is necessary that either  $f_W(0) = 0$  or  $f_Y(0) = 0$ .

Moreover, we want to have  $Y = \phi(X)$ , where  $X$  is Gaussian. In order to obtain  $Y$  with a zero density at 0, it is necessary to build a function  $\phi$  with  $\phi'(0) = \infty$  (see Appendix B), which is usually not desirable as an activation function of a neural network for training stability reasons. So, it is natural to build  $W$  such that  $f_W(0) = 0$ .

**Remark 1.** *The density of the product of two i.i.d.  $\mathcal{N}(0, 1)$  random variables is:*

$$f(x) = \frac{K_0(|x|)}{\pi},$$

where  $K_0$  is the modified Bessel function of the second kind. Naturally,  $\lim_{x \rightarrow 0} f(x) = \infty$ , which illustrates Corollary 1.

### 2.2.2 Behavior of the tail

We use the results of [8] on the “generalized Weibull-tail distributions”.

**Definition 1** (slowly varying function). *A measurable function  $l : (0, \infty) \rightarrow (0, \infty)$  is said to be “slowly varying” if:*

$$\forall a > 0, \quad \lim_{x \rightarrow \infty} \frac{l(ax)}{l(x)} = 1.$$

**Definition 2** (generalized Weibull-tail distribution). *A random variable  $X$  is called “generalized Weibull-tail” with tail parameter  $\theta > 0$  if its survival function  $S$  is bounded by Weibull-tail functions of tail parameter  $\theta$  with possibly different slowly-varying functions  $l_1$  and  $l_2$ :*

$$\exp(-x^\theta l_1(x)) \leq S(x) \leq \exp(-x^\theta l_2(x)), \quad \text{for all } x > 0.$$

**Proposition 2** (behavior of the tail). *The product of two independent non-negative generalized Weibull-tail random variables  $|W|$  and  $|Y|$  with tail parameters  $\theta_W$  and  $\theta_Y$  is generalized Weibull-tail with tail parameter  $\theta$  such that:*

$$\frac{1}{\theta} = \frac{1}{\theta_W} + \frac{1}{\theta_Y}.$$

We recall that, in our case,  $|G| = |W| \cdot |Y|$  is the absolute value of a Gaussian random variable. So  $|G|$  is generalized Weibull-tail of parameter  $\theta_G = 2$ . Thus, if we assume that  $|W|$  and  $|Y|$  are generalized Weibull-tail of respective parameters  $\theta_W$  and  $\theta_Y$ , then we necessarily have:

$$\frac{1}{\theta_Y} = \frac{1}{2} - \frac{1}{\theta_W}.$$

So  $\theta := \theta_W$  is constrained to the set  $\Theta = (2, \infty)$ .

## 2.3 Product of two random variables

We assume that  $G = WY \sim \mathcal{N}(0, 1)$ . Let us consider the random variables  $|W|$ ,  $|Y|$  and  $|G| = |W| \cdot |Y|$ . Let  $f_{|W|}$ ,  $f_{|Y|}$  and  $f_{|G|}$  their densities. Under integrability conditions, we can use the following property of the Mellin transform  $\mathcal{M}$ :

$$\mathcal{M}f_{|G|} = (\mathcal{M}f_{|W|}) \cdot (\mathcal{M}f_{|Y|}), \quad \text{where} \quad (\mathcal{M}f)(t) = \int_0^\infty x^{t-1} f(x) dx.$$

So:

$$f_{|Y|}(y) := \mathcal{M}^{-1} \left[ \frac{\mathcal{M}f_{|G|}}{\mathcal{M}f_{|W|}} \right] (y).$$

Then, by symmetry, we can obtain  $f_Y$  starting from  $f_{|G|}$  (density of the absolute value of a Gaussian  $\mathcal{N}(0, 1)$ ) and  $f_{|W|}$ . We choose  $W \sim \mathcal{W}(\theta, 1)$ , the symmetric Weibull distribution of parameters  $(\theta, 1)$  (see Eqn. (3)). To summarize:

$$W \sim P_\theta = \mathcal{W}(\theta, 1), \quad Y \sim Q_\theta : f_{|Y|} = \mathcal{M}^{-1} \left[ \frac{\mathcal{M}f_{|G|}}{\mathcal{M}f_{|W|}} \right] (y).$$

### 2.3.1 Obtaining the distribution of $Y$

**Through the Mellin transform.** The most direct way to obtain the distribution of  $|Y|$  consists in performing the following computation:

$$f_{|Y|}(y) := \mathcal{M}^{-1} \left[ \frac{\mathcal{M}f_{|G|}}{\mathcal{M}f_{|W|}} \right] (y).$$

Sadly, while  $\mathcal{M}f_{|G|}$  and  $\mathcal{M}f_{|W|}$  are easy to compute, the inverse Mellin transform  $\mathcal{M}^{-1}$  seems to be analytically untractable in this case.

Moreover, the numerical computation of  $\mathcal{M}^{-1}$  is too heavy for the required precision.

**Through the numerical resolution of an integral equation.** Provided that  $|G| = |W| \cdot |Y|$ , where  $G \sim \mathcal{N}(0, 1)$  and  $Y \sim \mathcal{W}(\theta, 1)$ , we are looking for the distribution of  $|Y|$ . We are able to express this problem as an integral equation, namely:

$$S_{|G|}(z) = \int_0^\infty f_{|Y|}(y) S_{|W|} \left( \frac{z}{y} \right) dy, \quad (7)$$

where  $f_{|Y|}$  is the unknown, and  $S_{|G|}$  and  $S_{|W|}$  are respectively the survival function of  $|G|$  and  $|W|$ .

We provide in Appendix D the technical details of the numerical resolution, based on the resolution of a least-square problem. However, we need to express the candidate solution as a linear combination of a family of functions. We propose the following family:

$$\begin{aligned} e_{\lambda_1}^{(1)}(x) &= \exp \left[ - \left( \frac{x}{\lambda_1} \right)^{\theta'} \right], \\ e_{\lambda_2, \mu_2, \alpha_2, \gamma_2}^{(2)}(x) &= \exp \left[ - \left( \frac{|x - \mu_2|}{\lambda_2} \right)^{\theta'} \right] (x + \gamma_2)^{\alpha_2}, \\ e_{\lambda_3, \mu_3, \alpha_3, \gamma_3}^{(3)}(x) &= \exp \left[ - \left( \frac{|x - \mu_3|}{\lambda_3} \right)^{\theta'} \right] (x + \gamma_3)^{-\alpha_3}, \\ e_{\lambda_4, \mu_4}^{(4)}(x) &= \exp \left[ - \left( \frac{|x - \mu_4|}{\lambda_4} \right)^{\theta'} \right], \end{aligned}$$

where  $\theta' = [\frac{1}{2} - \frac{1}{\theta}]^{-1}$  (see Proposition 2).

The resolution process consists in the following loop:

1. solve the integral equation with the family of functions:

$$(e_{\lambda_1}^{(1)}, e_{\lambda_2, \mu_2, \alpha_2, \gamma_2}^{(2)}, e_{\lambda_3, \mu_3, \alpha_3, \gamma_3}^{(3)}, e_{\lambda_4, \mu_4}^{(4)});$$

2. train the  $\lambda_i, \mu_i, \alpha_i, \gamma_i$  by gradient descent.

## 2.4 Obtaining the activation function

To get  $\phi$ , it is sufficient to make the following computation:

$$\phi_\theta(t) := F_Y^{-1}(F_G(t)), \quad F_Y(t) := \frac{1}{2} + \frac{1}{2} \operatorname{sgn}(t) \int_0^{|t|} f_{|Y|}(y) dy. \quad (8)$$

## 2.5 Results

**Examples of activation functions.** Using these activation functions  $\phi_\theta$ , combined with the adequate initialization of weights  $\mathcal{W}(\theta, 1)$ , guarantees Gaussian pre-activations. The family of the  $\phi_\theta$  is a continuum spanning bounded functions, such as Heaviside, or tanh, and linear functions.

## 3 Discussion

**Independence of the components  $Z_i^{(l)}$  of the pre-activations.** In this work, we assumed that the components  $Z_i^{(l)}$  of the pre-activations are independent. This simplification is partly justified by the result obtained in [4]: in this variant of the Central Limit Theorem, exchangeability of the components  $Z_i^{(l)}$  is sufficient to ensure that the pre-activations of the next layer  $Z^{(l+1)}$  are Gaussian. However, if we stay in a non-asymptotic case, where the width  $n_l$  of the layers is small, this result does not apply, and it is necessary to measure the effect of the dependence between the components of  $Z^{(l)}$ .

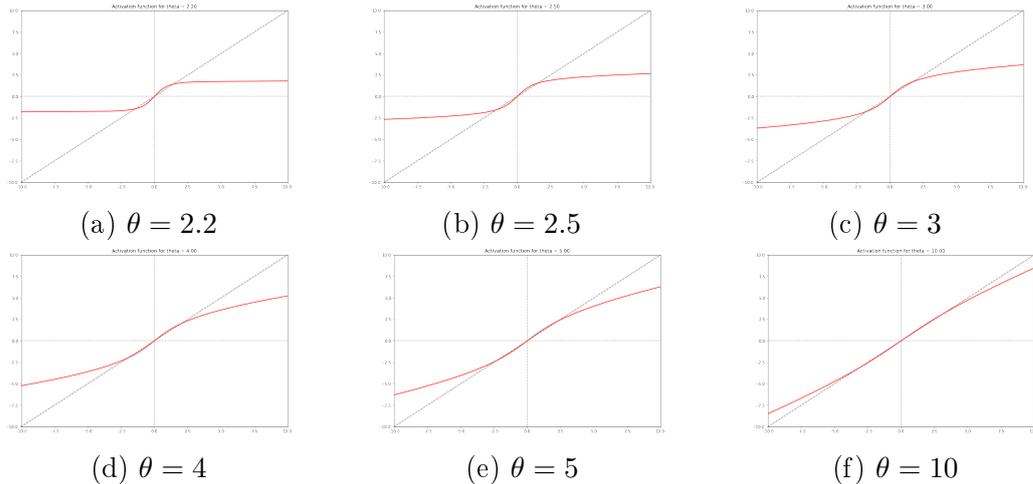


Figure 1 – Activation functions  $\phi_\theta$  with  $\theta \in \{2.2; 2.5; 3; 4; 5; 10\}$ .

**Initialization of the biases.** We have not taken into account the initialization of the biases, since they are just added to a sum, which is already Gaussian (see Eqn. (4)). We can set them to 0 or draw them randomly from a Gaussian, provided that we adjust the first term of Eqn. (4) in such a way that the result remains  $\mathcal{N}(0, 1)$ .

**Back to the Neural Tangent Kernels and the Edge of Chaos.** In these works, the hypothesis of Gaussian pre-activations come from the infinite-width limit of the layers. It is now possible to do the same work, but with the pairs initialization/activation we are proposing here, for narrow neural networks, while keeping the Gaussian hypothesis (necessary when deriving the results).

**Testing neural networks with several pairs  $(P_\theta, \phi_\theta)$ .** We are now ready to test the Gaussian hypothesis for the pre-activations by training a family of neural networks having each their own pair  $(P_\theta, \phi_\theta)$ . We would be able to compare their performance with networks with the same architecture, but with the usual initialization and activation functions. The results will help us to answer the current hot question: is it desirable to have Gaussian pre-activations in a neural network?

## Acknowledgements

The project leading to this work has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement No 834175).

## References

- [1] Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. *Advances in Neural Information Processing Systems*, 32, 2019.
- [2] Vincent Fortuin, Adrià Garriga-Alonso, Florian Wenzel, Gunnar Rätsch, Richard Turner, Mark van der Wilk, and Laurence Aitchison. Bayesian neural network priors revisited. *arXiv preprint arXiv:2102.06571*, 2021.
- [3] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in Neural Information Processing Systems*, 31, 2018.
- [4] Alexander G de G Matthews, Jiri Hron, Mark Rowland, Richard E Turner, and Zoubin Ghahramani. Gaussian process behaviour in wide deep neural networks. In *International Conference on Learning Representations*, 2018.
- [5] Ben Poole, Subhaneil Lahiri, Maithra Raghu, Jascha Sohl-Dickstein, and Surya Ganguli. Exponential expressivity in deep neural networks through transient chaos. *Advances in Neural Information Processing Systems*, 29, 2016.
- [6] Hadi Pouransari, Zhucheng Tu, and Oncel Tuzel. Least squares binary quantization of neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 698–699, 2020.
- [7] Frank Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological Review*, 65(6):386, 1958.
- [8] Mariia Vladimirova, Julyan Arbel, and Stéphane Girard. Bayesian neural network unit priors and generalized weibull-tail property. In *Asian Conference on Machine Learning*, pages 1397–1412. PMLR, 2021.
- [9] Mariia Vladimirova, Stéphane Girard, Hien Nguyen, and Julyan Arbel. Sub-Weibull distributions: Generalizing sub-Gaussian and sub-exponential properties to heavier tailed distributions. *Stat*, 9(1):e318, 2020.

## A Decomposition of a Gaussian as a sum of i.i.d. random variables

**Lemma 1.** *Let  $(Z_i)_i$  be i.i.d. random variables. Let  $Z = \frac{1}{\sqrt{n}} \sum_{i=1}^n Z_i$ . If  $Z$  is  $\mathcal{N}(0, 1)$ , then the distribution of the  $Z_i$  is also  $\mathcal{N}(0, 1)$ .*

*Proof.* Let  $\psi_Z(x) := \mathbb{E}[e^{iZx}]$  be the characteristic function of the distribution of  $Z$ . Since the  $(Z_i)_i$  are i.i.d., we have:

$$\begin{aligned}\psi_Z(x) &= \exp\left(-\frac{x^2}{2}\right), \\ \psi_Z(x) &= \left[\psi_{\frac{Z_i}{\sqrt{n}}}(x)\right]^n.\end{aligned}$$

Therefore,  $\psi_{Z_i}(x) = e^{-x^2/2}$ . Thus,  $Z_i \sim \mathcal{N}(0, 1)$ . □

## B Activation functions with vertical tangent at 0

**Lemma 2.** *Let  $\phi_\theta$  a function transforming a Gaussian random variable  $G \sim \mathcal{N}(0, 1)$  into a symmetrical Weibull random variable  $Y \sim \mathcal{W}(\theta, 1)$ . That is,  $Y = \phi_\theta(G)$ . Then  $\phi_\theta$  has a vertical tangent at 0.*

*Proof.* We have:

$$\phi_\theta(x) = F_Y^{-1}(F_G(x)),$$

where:

$$\begin{aligned}F_G(x) &:= \frac{2}{\sqrt{2\pi}} \int_{-\infty}^x \exp\left(-\frac{t^2}{2}\right) dt \\ F_Y(x) &:= \frac{1}{2} + \frac{1}{2} \operatorname{sgn}(x) \exp(-|x|^\theta).\end{aligned}$$

Thus:

$$\phi'_\theta(x) = F'_G(x) \frac{1}{F'_Y(F_Y^{-1}(F_G(x)))}$$

Therefore:

$$\begin{aligned}\phi'_\theta(0) &= F'_G(0) \frac{1}{F'_Y(F_Y^{-1}(F_G(0)))} \\ &= \frac{1}{\sqrt{2\pi}} \frac{1}{F'_Y(F_Y^{-1}(1/2))} \\ &= \frac{1}{\sqrt{2\pi}} \frac{1}{F'_Y(0)} \\ &= \infty.\end{aligned}$$

□

## C Constraints on the product of two random variables

**Proposition 1** (density of a product of random variables at 0). *Let  $W, Y$  be two independent non-negative random variables and  $Z = WY$ . Let  $f_W, f_Y, f_Z$  be their respective density. Assuming that  $f_Y$  is continuous at 0 with  $f_Y(0) > 0$ , we have:*

$$\lim_{w \rightarrow 0} \int_w^\infty \frac{f_W(t)}{t} dt = \infty \quad \Rightarrow \quad \lim_{z \rightarrow 0} f_Z(z) = \infty, \quad (9)$$

moreover, if  $f_Y$  is bounded:

$$\int_0^\infty \frac{f_W(t)}{t} dt < \infty \quad \Rightarrow \quad f_Z(0) = f_Y(0) \int_0^\infty \frac{f_W(t)}{t} dt. \quad (10)$$

*Proof.* Let  $z, z_0 > 0$ :

$$\begin{aligned} f_Z(z) &= \int_0^\infty f_Y(t) \frac{1}{t} f_W\left(\frac{z}{t}\right) dt \\ &\geq \int_0^{z_0} f_Y(t) \frac{1}{t} f_W\left(\frac{z}{t}\right) dt \\ &\geq \inf_{[0, z_0]} f_Y \cdot \int_0^{z_0} \frac{1}{t} f_W\left(\frac{z}{t}\right) dt \\ &\geq \inf_{[0, z_0]} f_Y \cdot \int_{z/z_0}^\infty \frac{f_W(t)}{t} dt. \end{aligned}$$

Let us take  $z_0 = \sqrt{z}$ . We have:

$$f_Z(z) \geq \inf_{[0, \sqrt{z}]} f_Y \cdot \int_{\sqrt{z}}^\infty \frac{f_W(t)}{t} dt.$$

Then we take the limit  $z \rightarrow 0$ , hence:

- if  $\int_0^\infty \frac{f_W(t)}{t} dt = \infty$ , then:  $\lim_{z \rightarrow 0} f_Z(z) = \infty$ , which achieves (9);
- if  $\int_0^\infty \frac{f_W(t)}{t} dt < \infty$ , then:  $f_Z(0) \geq f_Y(0) \int_0^\infty \frac{f_W(t)}{t} dt$ , which achieves one half of (10);

Let us prove the second half of (10). Let  $z, z_0 > 0$ :

$$\begin{aligned} f_Z(z) &= \int_0^\infty f_Y(t) \frac{1}{t} f_W\left(\frac{z}{t}\right) dt \\ &= \int_0^{z_0} f_Y(t) \frac{1}{t} f_W\left(\frac{z}{t}\right) dt + \int_{z_0}^\infty f_Y(t) \frac{1}{t} f_W\left(\frac{z}{t}\right) dt \\ &\leq \sup_{[0, z_0]} f_Y \cdot \int_{z/z_0}^\infty \frac{f_W(t)}{t} dt + \int_1^\infty f_Y(z_0 t) \frac{1}{t} f_W\left(\frac{z}{z_0 t}\right) dt \end{aligned}$$

Let  $z_0 = \sqrt{z}$ . We have:

$$f_Z(z) \leq \sup_{[0, \sqrt{z}]} f_Y \cdot \int_{\sqrt{z}}^\infty \frac{f_W(t)}{t} dt + \int_1^\infty f_Y(\sqrt{z}t) \frac{1}{t} f_W\left(\frac{\sqrt{z}}{t}\right) dt,$$

where:

$$\begin{aligned} \int_1^\infty f_Y(\sqrt{zt}) \frac{1}{t} f_W\left(\frac{\sqrt{z}}{t}\right) dt &\leq \|f_Y\|_\infty \int_1^\infty \frac{1}{t} f_W\left(\frac{\sqrt{z}}{t}\right) dt \\ &\leq \|f_Y\|_\infty \int_0^{\sqrt{z}} \frac{f_W(t)}{t} dt. \end{aligned}$$

According to the hypotheses, we have, as  $z \rightarrow 0$ :

$$\begin{aligned} \sup_{[0, \sqrt{z}]} f_Y &\rightarrow f_Y(0) \\ \int_{\sqrt{z}}^\infty \frac{f_W(t)}{t} dt &\rightarrow \int_0^\infty \frac{f_W(t)}{t} dt \\ \int_0^{\sqrt{z}} \frac{f_W(t)}{t} dt &\rightarrow 0, \end{aligned}$$

hence the result. □

## D Numerical resolution of an integral equation

We consider the following equation, of unknown  $g$ :

$$f(z) = \int_0^\infty g(x) K(z, x) dx, \quad (11)$$

where  $f$  is a fixed function and  $K$  is a fixed kernel (e.g.,  $f = S_Z$  and  $K(z, x) = S(z/x)$ ).

Let  $(z_1, \dots, z_n)$  a finite sequence of points. Let  $(e_1, \dots, e_p)$  a family of functions. We are looking for  $\hat{g} \in \mathcal{V}_p = \text{Vect}(e_1, \dots, e_p)$  such that:

$$\begin{aligned} \hat{g} &:= \sum_{j=1}^p c_j e_j, \\ \hat{g} &= \arg \min_{g \in \mathcal{V}_p} \sum_{i=1}^p \left( \frac{f(z_i) - \int_0^\infty g(x) K(z_i, x) dx}{f(z_i)} \right)^2, \end{aligned}$$

if  $f$  is strictly positive on  $\mathbb{R}^+$ .

This problem is a standard weighted least-square regression with solution:

$$\begin{aligned} C &= (M^T \Delta^2 M)^{-1} M^T \Delta^2 F \\ &= (M'^T M')^{-1} M'^T F', \end{aligned}$$

where:

$$\begin{array}{ll} C \in \mathbb{R}^p & C_j = c_j, \\ M \in \mathbb{R}^{n \times p} & M_{ij} = \int_0^\infty e_j(x) K(z_i, x) dx, \\ F \in \mathbb{R}^n & F_i = f(z_i), \\ \Delta \in \mathbb{R}^{n \times n} & \Delta = \text{Diag}(f(z_1)^{-1}, \dots, f(z_n)^{-1}), \\ M' = \Delta M & F' = \Delta F. \end{array}$$

This gives the approximate solution  $\hat{g}$  of Eqn. (11) minimizing the squared relative error between  $f$  and  $\hat{f}(\cdot) = \int \hat{g}(x)K(\cdot, x) dx$ .