



Understanding Priors in Bayesian Neural Networks at the Unit Level

Mariia Vladimirova, Jakob Verbeek, Pablo Mesejo, Julyan Arbel

► To cite this version:

Mariia Vladimirova, Jakob Verbeek, Pablo Mesejo, Julyan Arbel. Understanding Priors in Bayesian Neural Networks at the Unit Level. ICML 2019 - 36th International Conference on Machine Learning, Jun 2019, Long Beach, United States. pp.6458-6467. hal-02177151

HAL Id: hal-02177151

<https://hal.science/hal-02177151>

Submitted on 8 Jul 2019

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.

Understanding Priors in Bayesian Neural Networks at the Unit Level

Mariia Vladimirova^{1,2} Jakob Verbeek¹ Pablo Mesejo³ Julyan Arbel¹

Abstract

We investigate deep Bayesian neural networks with Gaussian weight priors and a class of ReLU-like nonlinearities. Bayesian neural networks with Gaussian priors are well known to induce an \mathcal{L}^2 , “weight decay”, regularization. Our results characterize a more intricate regularization effect at the level of the unit activations. Our main result establishes that the induced prior distribution on the units before and after activation becomes increasingly heavy-tailed with the depth of the layer. We show that first layer units are Gaussian, second layer units are sub-exponential, and units in deeper layers are characterized by sub-Weibull distributions. Our results provide new theoretical insight on deep Bayesian neural networks, which we corroborate with experimental simulation results.

1. Introduction

Neural networks (NNs), and their deep counterparts (Goodfellow et al., 2016), have largely been used in many research areas such as image analysis (Krizhevsky et al., 2012), signal processing (Graves et al., 2013), or reinforcement learning (Silver et al., 2016), just to name a few. The impressive performance provided by such machine learning approaches has greatly motivated research that aims at a better understanding the driving mechanisms behind their effectiveness. In particular, the study of the NNs distributional properties through Bayesian analysis has recently gained much attention.

Bayesian approaches investigate models by assuming a prior distribution on their parameters. Bayesian machine learning

refers to extending standard machine learning approaches with posterior inference, a line of research pioneered by works on Bayesian neural networks (Neal, 1992; MacKay, 1992). There is a large variety of applications, e.g. gene selection (Liang et al., 2018), and the range of models is now very broad, including e.g. Bayesian generative adversarial networks (Saatci & Wilson, 2017). See Polson & Sokolov (2017) for a review. The interest of the Bayesian approach to NNs is at least twofold. First, it offers a principled approach for modeling uncertainty of the training procedure, which is a limitation of standard NNs which only provide point estimates. A second main asset of Bayesian models is that they represent regularized versions of their classical counterparts. For instance, maximum a posteriori (MAP) estimation of a Bayesian regression model with double exponential (Laplace) prior is equivalent to Lasso regression (Tibshirani, 1996), while a Gaussian prior leads to ridge regression. When it comes to NNs, the regularization mechanism is also well appreciated in the literature, since they traditionally suffer from overparameterization, resulting in overfitting.

Central in the field of regularization techniques is the *weight decay* penalty (Krogh & Hertz, 1991), which is equivalent to MAP estimation of a Bayesian neural network with independent Gaussian priors on the weights. Dropout has recently been suggested as a regularization method in which neurons are randomly turned off (Srivastava et al., 2014), and Gal & Ghahramani (2016) proved that a neural network trained with dropout is equivalent to a probabilistic model, i.e. a deep Gaussian process (Damianou & Lawrence, 2013), leading to the consideration of such NNs as Bayesian models.

This paper is devoted to the investigation of hidden units prior distributions in Bayesian neural networks under the assumption of independent Gaussian weights. We first describe a fully connected neural network architecture as illustrated in Figure 1. Given an input $\mathbf{x} \in \mathbb{R}^N$, the ℓ -th hidden layer unit activations are defined as

$$\mathbf{g}^{(\ell)}(\mathbf{x}) = \mathbf{W}^{(\ell)} \mathbf{h}^{(\ell-1)}(\mathbf{x}), \quad \mathbf{h}^{(\ell)}(\mathbf{x}) = \phi(\mathbf{g}^{(\ell)}(\mathbf{x})), \quad (1)$$

where $\mathbf{W}^{(\ell)}$ is a weight matrix including the bias vector. A nonlinear activation function $\phi : \mathbb{R} \rightarrow \mathbb{R}$ is applied element-wise, which is called nonlinearity, $\mathbf{g}^{(\ell)} = \mathbf{g}^{(\ell)}(\mathbf{x})$ is a vector of pre-nonlinearities, and $\mathbf{h}^{(\ell)} = \mathbf{h}^{(\ell)}(\mathbf{x})$ is a

¹Univ. Grenoble Alpes, Inria, CNRS, Grenoble INP, LJK, 38000 Grenoble, France.

²Moscow Institute of Physics and Technology, 141701 Dolgoprudny, Russia.

³Andalusian Research Institute in Data Science and Computational Intelligence (DaSCI), University of Granada, 18071 Granada, Spain.

Correspondence to: Mariia Vladimirova <mariia.vladimirova@inria.fr>.

vector of post-nonlinearities. When we refer to either pre- or post-nonlinearities, we will use the notation $U^{(\ell)}$.

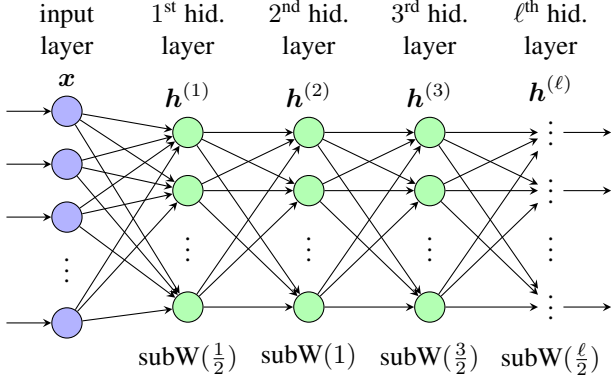


Figure 1. Neural network architecture and characterization of the ℓ -layer units prior distribution as sub-Weibull distribution with tail parameter $\ell/2$, see Definition 2.2.

1.1. Contributions

In this paper, we extend the theoretical understanding of feedforward fully connected NNs by studying prior distributions at the units level, under the assumption of independent and normally distributed weights. Our contributions are the following:

- (i) As our main contribution, we prove in Theorem 2.1 that under some conditions on the activation function ϕ , a Gaussian prior on the weights induces a sub-Weibull distribution on the units (both pre- and post-nonlinearities) with optimal tail parameter $\theta = \ell/2$, see Figure 1. The condition on ϕ essentially imposes that ϕ strikes at a linear rate to $+\infty$ or $-\infty$ for large absolute values of the argument, as ReLU does. In the case of bounded support ϕ , like sigmoid or tanh, the units are bounded, making them *de facto* sub-Gaussian⁴.
- (ii) We offer an interpretation of the main result from a more elaborate regularization scheme at the level of the units in Section 3.

1.2. Related work

Studying the distributional behaviour of feedforward networks has been a fruitful avenue for understanding these models, as pioneered by the works of Radford Neal (Neal, 1992; 1996) and David MacKay (MacKay, 1992). The first results in the field addressed the limiting setting when the number of units per layer tends to infinity, also called the wide regime. Neal (1996) proved that a single hidden layer neural network with normally distributed weights

tends in distribution in the wide limit either to a Gaussian process (Rasmussen & Williams, 2006) or to an α -stable process, depending on how the prior variance on the weights is rescaled. In recent works, Matthews et al. (2018b), or its extended version Matthews et al. (2018a), and Lee et al. (2018) extend the result of Neal to more-than-one-layer neural networks: when the number of hidden units grows to infinity, deep neural networks (DNNs) also tend in distribution to the Gaussian process, under the assumption of Gaussian weights for properly rescaled prior variances. For the rectified linear unit (ReLU) activation function, the Gaussian process covariance function is obtained analytically (Cho & Saul, 2009). For other nonlinear activation functions, Lee et al. (2018) use a numerical approximation algorithm.

Various distributional properties are also studied in NNs regularization methods. The *dropout* technique (Srivastava et al., 2014) was reinterpreted as a form of approximate Bayesian variational inference (Kingma et al., 2015; Gal & Ghahramani, 2016). While Gal & Ghahramani (2016) built a connection between dropout and the Gaussian process, Kingma et al. (2015) proposed a way to interpret Gaussian dropout. They suggested *variational dropout* where each weight of a model has its individual dropout rate. *Sparse variational dropout* (Molchanov et al., 2017) extends variational dropout to all possible values of dropout rates, and leads to a sparse solution. The approximate posterior is chosen to factorize either over rows or individual entries of the weight matrices. The prior usually factorizes in the same way, and the choice of the prior and its interaction with the approximating posterior family are studied by Hron et al. (2018). Performing dropout can be used as a Bayesian approximation but, as noted by Duvenaud et al. (2014), it has no regularization effect on infinitely-wide hidden layers.

Recent work by Bibi et al. (2018) provides the expression of the first two moments of the output units of a one layer NN. Obtaining the moments is a first step towards characterizing the full distribution. However, the methodology of Bibi et al. (2018) is limited to the first two moments and to single-layer NNs, while we address the problem in more generality for deep NNs.

In the remainder of the paper, we present our main contributions starting with the necessary statistical background and theoretical results (i), then moving to intuitions and interpretation (ii), and ending up with the description of the experiments and the discussion of the results obtained. More specifically, Section 2 states our main contribution, Theorem 2.1, with a proof sketch while additional technical results are deferred to Supplementary material. Section 3 illustrates penalization techniques, providing an interpretation for the theorem. Section 4 describes the experiments. Conclusions and directions for future work are presented in

⁴A trivial version of our main result holds, see Remark 2.1.

Section 5.

2. Bayesian neural networks have heavy-tailed deep units

The deep learning approach uses stochastic gradient descent and error back-propagation in order to fit the network parameters $(\mathbf{W}^{(\ell)})_{1 \leq \ell \leq L}$, where ℓ iterates over all network layers. In the Bayesian approach, the parameters are random variables described by probability distributions.

2.1. Assumptions on neural network

We assume a prior distribution on the model parameters, that are the weights \mathbf{W} . In particular, let all weights (including biases) be independent and have zero-mean normal distribution

$$W_{i,j}^{(\ell)} \sim \mathcal{N}(0, \sigma_w^2), \quad (2)$$

for all $1 \leq \ell \leq L$, $1 \leq i \leq H_{\ell-1}$ and $1 \leq j \leq H_\ell$, with fixed variance σ_w^2 . Given some input \mathbf{x} , such prior distribution induces by forward propagation (1) a prior distribution on the pre-nonlinearity and post-nonlinearity, whose *tail properties* are the focus of this section. To this aim, the nonlinearity ϕ is required to span at least half of the real line as follows. We introduce an extended version of the nonlinearity assumption from Matthews et al. (2018a):

Definition 2.1 (Extended envelope property for nonlinearities). *A nonlinearity $\phi : \mathbb{R} \rightarrow \mathbb{R}$ is said to obey the extended envelope property if there exist $c_1, c_2 \geq 0$, $d_1, d_2 > 0$ such that the following inequalities hold*

$$\begin{aligned} |\phi(u)| &\geq c_1 + d_1|u| \quad \text{for all } u \in \mathbb{R}_+ \text{ or } u \in \mathbb{R}_-, \\ |\phi(u)| &\leq c_2 + d_2|u| \quad \text{for all } u \in \mathbb{R}. \end{aligned} \quad (3)$$

The interpretation of this property is that ϕ must shoot to infinity at least in one direction (\mathbb{R}_+ or \mathbb{R}_- , at least linearly (first line of (3)), and also at most linearly (second line of (3)). Of course, compactly supported nonlinearities such as sigmoid and tanh do not satisfy the extended envelope property but the majority of other nonlinearities do, including ReLU, ELU, PReLU, and SeLU.

Lemma 2.1. *Let a nonlinearity $\phi : \mathbb{R} \rightarrow \mathbb{R}$ obey the extended envelope property. Then for any symmetric random variable X the following asymptotic equivalence⁵ holds*

$$\|\phi(X)\|_k \asymp \|X\|_k, \quad \text{for all } k \geq 1, \quad (4)$$

where $\|X\|_k = (\mathbb{E}[|X|^k])^{1/k}$ is a k -th norm of X .

⁵See Definition B.1 for the asymptotic equivalence \asymp definition in Supplementary material.

Proof. According to asymptotic equivalence definition there must exist positive constants d and D such that for all $k \in \mathbb{N}$ it holds

$$d \leq \|\phi(X)\|_k / \|X\|_k \leq D. \quad (5)$$

The extended envelope property upper bound and the triangle inequality for norms imply the right-hand side of (5), since

$$\|\phi(X)\|_k \leq \|c_2 + d_2|u|\|_k \leq c_2 + d_2\|u\|_k.$$

Assume that $|\phi(u)| \geq c_1 + d_1|u|$ for $u \in \mathbb{R}_+$. Consider the lower bound of the nonlinearity moments

$$\|\phi(X)\|_k \geq \|d_1 u_+\|_k + c_1 + \|\phi(u_-)\|_k,$$

where $\{u_- : u \in \mathbb{R}_-\}$ and $\{u_+ : u \in \mathbb{R}_+\}$. For negative u_- there are constants $c_1 \geq 0$ and $d_1 > 0$ such that $c_1 - d_1 u > |\phi(u)|$, or $c_1 > |\phi(u)| + d_1 u$:

$$\|\phi(X)\|_k > \|d_1 u_+\|_k + \|\phi(u_-)\|_k + d_1 u_- \|u\|_k \geq d_1 \|u\|_k.$$

It yields asymptotic equivalence (4). \square

2.2. Main theorem

This section postulates the rigorous result with a proof sketch. In Supplementary material one can find proofs of intermediate lemmas.

Firstly, we define the notion of *sub-Weibull* random variables.

Definition 2.2 (Sub-Weibull random variable). *A random variable X satisfying for all $x > 0$ and for some $\theta > 0$*

$$\mathbb{P}(|X| \geq x) \leq a \exp(-x^{1/\theta}), \quad (6)$$

is called a sub-Weibull random variable with so-called tail parameter θ , which is denoted by $X \sim \text{subW}(\theta)$.

Sub-Weibull distributions are characterized by tails lighter than (or equally light as) Weibull distributions; in the same way as sub-Gaussian or sub-exponential distributions correspond to distributions with tails lighter than Gaussian and exponential distributions, respectively. Sub-Weibull distributions are parameterized by a positive tail index θ and are equivalent to sub-Gaussian for $\theta = 1/2$ and sub-exponential for $\theta = 1$. More information on sub-Weibull distributions can be found in Supplementary material. To describe a tail lower bound through some sub-Weibull distribution family, i.e. a distribution of X to have the tail heavier than some sub-Weibull, we define an optimal tail parameter for that distribution as follows:

Proposition 2.1 (Optimal sub-Weibull tail coefficient and moment condition). *Let $\theta > 0$ and let X be a random variable satisfying the following asymptotic equivalence on moments*

$$\|X\|_k \asymp k^\theta, \quad \text{for all } k \geq 1.$$

Then X is sub-Weibull distributed with optimal tail parameter θ , in the sense that for any $\theta' < \theta$, X is not sub-Weibull with tail parameter θ' .

The following theorem postulates the main results.

Theorem 2.1 (Sub-Weibull units). *Consider a feed-forward Bayesian neural network with Gaussian priors (2) and with nonlinearity ϕ satisfying the extended envelope condition of Definition 2.1. Then conditional on the input \mathbf{x} , the marginal prior distribution induced by forward propagation (1) on any unit (pre- or post-nonlinearity) of the ℓ -th hidden layer is sub-Weibull with optimal tail parameter $\theta = \ell/2$. That is for any $1 \leq \ell \leq L$, and for any $1 \leq m \leq H_\ell$,*

$$U_m^{(\ell)} \sim \text{subW}(\ell/2),$$

where a **subW** distribution is defined in Definition 2.2, and $U_m^{(\ell)}$ is either a pre-nonlinearity $g_m^{(\ell)}$ or a post-nonlinearity $h_m^{(\ell)}$.

Proof. The idea is to prove by induction with respect to hidden layer depth ℓ that pre- and post-nonlinearity satisfy the asymptotic moment equivalence

$$\|g^{(\ell)}\|_k \asymp k^{\ell/2} \text{ and } \|h^{(\ell)}\|_k \asymp k^{\ell/2}.$$

The statement of the theorem then follows by the moment characterization of optimal sub-Weibull tail coefficient in Proposition 2.1.

According to Lemma A.1 from Supplementary material, centering does not harm tail properties, then, for simplicity, we consider zero-mean distributions $W_{i,j}^{(\ell)} \sim \mathcal{N}(0, \sigma_w^2)$.

Base step: consider the distribution of the first hidden layer pre-nonlinearity g ($\ell = 1$). Since weights \mathbf{W}_m follow normal distribution and \mathbf{x} is a feature vector, then each hidden unit $\mathbf{W}_m^\top \mathbf{x}$ follow also normal distribution

$$g = \mathbf{W}_m^\top \mathbf{x} \sim \mathcal{N}(0, \sigma_w^2 \|\mathbf{x}\|^2).$$

Then, for normal zero-mean variable g , having variance $\sigma^2 = \sigma_w^2 \|\mathbf{x}\|^2$, holds the equality in sub-Gaussian property with variance proxy equals to normal distribution variance and from Lemma B.1 in Supplementary material:

$$\|g\|_k \asymp \sqrt{k}.$$

As activation function ϕ obeys extended envelope property, nonlinearity moments are asymptotic equivalent to symmetric variable moments

$$\|\phi(g)\|_k \asymp \|g\|_k \sim \sqrt{k}.$$

It implies that first hidden layer post-nonlinearity h have sub-Gaussian distribution or sub-Weibull with tail parameter $\theta = 1/2$ (Definition 2.2).

Inductive step: show that if the statement holds for $\ell - 1$, then it also holds for ℓ .

Suppose the post-nonlinearity of $(\ell - 1)$ -th hidden layer satisfies the moment condition. Hidden units satisfy the non-negative covariance theorem (Theorem 2.2):

$$\text{Cov} \left[\left(h^{(\ell-1)} \right)^s, \left(\tilde{h}^{(\ell-1)} \right)^t \right] \geq 0, \text{ for any } s, t \in \mathbb{N}.$$

Let the number of hidden units in $(\ell - 1)$ -th layer equals to H . Then according to Lemma B.2 from Supplementary material, under assumption of zero-mean Gaussian weights, pre-nonlinearity of ℓ -th hidden layer $g^{(\ell)} = \sum_{i=1}^H W_{m,i}^{(\ell-1)} h_i^{(\ell-1)}$ also satisfy the moment condition, but with $\theta = \ell/2$

$$\|g^{(\ell)}\|_k \asymp k^{\ell/2}.$$

Using the extended envelope property 2.1, one can show from (5) that post-nonlinearity $h^{(\ell)}$ satisfy the same moment condition as pre-nonlinearity $g^{(\ell)}$. This finishes the proof. \square

Remark 2.1. If the activation function ϕ is bounded, such as the sigmoid, or tanh, then the units are bounded. As a result, by Hoeffding's Lemma, they have a sub-Gaussian distribution.

Remark 2.2. Normalization techniques, such as batch normalization (Ioffe & Szegedy, 2015), layer normalization (Ba et al., 2016), significantly reduce the training time in feed-forward neural networks. It can be decomposed into elementary operations. According to Proposition A.4 from Supplementary material, elementary operations do not harm the distribution tail parameter. Therefore, normalization methods do not have an influence on tail behavior.

2.3. Intermediate theorem

This section states with a proof sketch that the covariance between hidden units in the neural network is non-negative.

Theorem 2.2 (Non-negative covariance between hidden units). *Consider the deep neural network described in, and with the assumptions of, Theorem 2.1. The covariance between hidden units of the same layer is non-negative. Moreover, for given ℓ -th hidden layer units $h^{(\ell)}$ and $\tilde{h}^{(\ell)}$, it holds*

$$\text{Cov} \left[\left(h^{(\ell)} \right)^s, \left(\tilde{h}^{(\ell)} \right)^t \right] \geq 0, \text{ where } s, t \in \mathbb{N}.$$

For first hidden layer $\ell = 1$ there is equality for all s and t .

Proof. A more detailed proof can be found in Supplementary material in Section 3.

Recall the covariance definition for random variables X and Y

$$\text{Cov}[X, Y] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]. \quad (7)$$

The proof is based on induction with respect to the hidden layer number.

In the proof let us make notation simplifications: $\mathbf{w}_m^\ell = \mathbf{W}_m^\ell$ and $w_{mi}^\ell = W_{mi}^\ell$ for all $m \in H_\ell$. If the index m is omitted, then \mathbf{w}^ℓ is some the vectors \mathbf{w}_m^ℓ , w_i^ℓ is i -th element of the vector \mathbf{w}_m^ℓ .

1. First hidden layer. Consider the first hidden layer units $h^{(1)}$ and $\tilde{h}^{(1)}$. The covariance between units is equal to zero and the units are Gaussian, since the weights $\mathbf{w}^{(1)}$ and $\tilde{\mathbf{w}}^{(1)}$ are from $\mathcal{N}(0, \sigma_w^2)$ and independent. Thus, the first hidden layer units are independent and its covariance (7) equals to 0. Moreover, since $h^{(1)}$ and $\tilde{h}^{(1)}$ are independent, then $(h^{(1)})^s$ and $(\tilde{h}^{(1)})^t$ are also independent.

2. Next hidden layers. Assume that $(\ell - 1)$ -th hidden layer has $H_{\ell-1}$ hidden units, where $\ell > 1$. Then ℓ -th hidden layer pre-nonlinearity is equal to

$$g^{(\ell)} = \sum_{i=1}^{H_{\ell-1}} w_i^{(\ell)} h_i^{(\ell-1)}. \quad (8)$$

We want to prove that the covariance (7) between ℓ -th hidden layer pre-nonlinearity is non-negative. Let us show firstly the idea of the proof in the case $H_{\ell-1} = 1$ and then briefly the proof for any finite $H_{\ell-1} > 1$, $H_{\ell-1} \in \mathbb{N}$.

2.1 One hidden unit. In the case $H_{\ell-1} = 1$, the covariance (7) sign is the same as of expression

$$\mathbb{E} \left[\left(h^{(\ell-1)} \right)^{2(s_1+t_1)} \right] - \mathbb{E} \left[\left(h^{(\ell-1)} \right)^{2s_1} \right] \mathbb{E} \left[\left(h^{(\ell-1)} \right)^{2t_1} \right],$$

since the weights are zero-mean distributed, its moments are equal to zero with odd order. According to Jensen's inequality for convex function f , we have $\mathbb{E}[f(x_1, x_2)] \geq f(\mathbb{E}[x_1], \mathbb{E}[x_2])$. Since a function $f(x_1, x_2) = x_1 x_2$ is convex for $x_1 \geq 0$ and $x_2 \geq 0$, then, taking $x_1 = (h^{(\ell-1)})^{2s_1}$ and $x_2 = (h^{(\ell-1)})^{2t_1}$, we have the condition we need (9) is satisfied.

2.1. H hidden units. Now let us consider the covariance between pre-nonlinearity (8) for $H_{\ell-1} = H > 1$. Raise the sum in the brackets to the power

$$\begin{aligned} \left(\sum_{i=1}^H w_i^{(\ell)} h_i^{(\ell-1)} \right)^s &= \\ &= \sum_{s_H=0}^s C_{s_H}^s \left(w_H^{(\ell)} h_H^{(\ell-1)} \right)^{s_H} \left(\sum_{i=1}^{H-1} w_i^{(\ell)} h_i^{(\ell-1)} \right)^{s-s_H}. \end{aligned}$$

And the same way for the second bracket And the same way for the second bracket $\left(\sum_{i=1}^H \tilde{w}_i^{(\ell)} h_i^{(\ell-1)} \right)^t$. Notice that binomial terms will be the same in the minuend and the

subtrahend terms of (7). So the covariance in our notations can be written in the form of

$$\begin{aligned} \text{Cov} \left[\left(\sum_{i=1}^{H_{\ell-1}} w_i^{(\ell)} h_i^{(\ell-1)} \right)^s, \left(\sum_{i=1}^{H_{\ell-1}} \tilde{w}_i^{(\ell)} h_i^{(\ell-1)} \right)^t \right] &= \\ &= \sum \sum C (\mathbb{E}[AB] - \mathbb{E}[A] \mathbb{E}[B]), \end{aligned}$$

where C -terms contain binomial coefficients, A -terms — all possible products of hidden units in $(g^{(\ell)})^s$ and B -terms — all possible products of hidden units in $(\tilde{g}^{(\ell)})^t$. For covariance being non-negative it is enough to show that the difference $\mathbb{E}[AB] - \mathbb{E}[A] \mathbb{E}[B]$ is non-negative. Since the weights are Gaussian and independent, we have the following equation, omitting the superscript for simplicity,

$$\begin{aligned} \mathbb{E}[AB] &= W \tilde{W} \cdot \mathbb{E} \left[\prod_{i=1}^H h_i^{s_i+t_i} \right], \\ \mathbb{E}[A] \mathbb{E}[B] &= W \tilde{W} \cdot \mathbb{E} \left[\prod_{i=1}^H h_i^{s_i} \right] \mathbb{E} \left[\prod_{i=1}^H h_i^{t_i} \right], \end{aligned}$$

where $W \tilde{W}$ is the product of weights moments

$$W \tilde{W} = \prod_{i=1}^H \mathbb{E}[w_i^{s_i}] \mathbb{E}[\tilde{w}_i^{t_i}].$$

For $W \tilde{W}$ not equal to zero, all the powers must be even. Now we need to prove

$$\mathbb{E} \left[\prod_{i=1}^{H/2} h_i^{2(s_i+t_i)} \right] \geq \mathbb{E} \left[\prod_{i=1}^{H/2} h_i^{2s_i} \right] \mathbb{E} \left[\prod_{i=1}^{H/2} h_i^{2t_i} \right] \quad (9)$$

According to Jensen's inequality for convex function, since a function $f(x_1, x_2) = x_1 x_2$ is convex for $x_1 \geq 0$ and $x_2 \geq 0$, then, taking $x_1 = \prod_{i=1}^{H/2} h_i^{2s_i}$ and $x_2 = \prod_{i=1}^{H/2} h_i^{2t_i}$, the condition from (9) is satisfied.

3. Post-nonlinearity.

Let show the proof for the ReLU nonlinearity.

The distribution of the ℓ -th hidden layer pre-nonlinearity $g^{(\ell)}$ is the sum of symmetric distributions, which are products of Gaussian variables $w^{(\ell)}$ and non-negative ReLU output, i.e. $(\ell - 1)$ -th hidden layer post-nonlinearity $h^{(\ell-1)}$. It leads that $g^{(\ell)}$ follows symmetric distribution and the following inequality

$$\begin{aligned} \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} g g' p(g, g') dg dg' &\geq \\ &\geq \int_{-\infty}^{+\infty} g p(g) dg \cdot \int_{-\infty}^{+\infty} g' p(g') dg' \end{aligned}$$

implies the same inequality for positive part

$$\begin{aligned} \int_0^{+\infty} \int_0^{+\infty} gg' p(g, g') dg dg' &\geq \\ &\geq \int_0^{+\infty} gp(g) dg \cdot \int_0^{+\infty} g' p(g') dg'. \end{aligned}$$

Notice that the equality above is the ReLU function output and for symmetric distribution we have

$$\int_0^{+\infty} xp(x) dx = \frac{1}{2} \mathbb{E}[|X|]. \quad (10)$$

That means if non-negative covariance is proven for pre-nonlinearities, for post-nonlinearities it is also non-negative. Omit the proof for the other nonlinearities with extended envelope property, since instead of precise equation (10), the asymptotic equivalence for moments will be used for positive part and for negative part — precise expectation expression which depends on certain nonlinearity. \square

2.4. Convolutional neural networks

Convolutional neural networks (Fukushima & Miyake, 1982; LeCun et al., 1998) are a particular kind of neural network for processing data that has a known grid-like topology, which allows to encode certain properties into the architecture. These then make the forward function more efficient to implement and vastly reduce the amount of parameters in the neural network. Neurons in such networks are arranged in three dimensions: width, height and depth. There are three main types of layers that can be concatenated in these architectures: convolutional, pooling, and fully-connected layers (exactly as seen in standard NNs). The convolutional layer computes dot products between a region in the inputs and its weights. Therefore, each region can be considered as a particular case of a fully-connected layer. Pooling layers control overfitting and computations in deep architectures. They operate independently on every slice of the input and reduces it spatially. The most commonly functions used in pooling layers are *max pooling* and *average pooling*.

Proposition 2.2. *The operations: 1. max pooling and 2. averaging do not modify the optimal tail parameter θ of sub-Weibull family. Consequently, the result of Theorem 2.1 carries over to convolutional neural networks.*

Proof. Let $X_i \sim \text{subW}(\theta)$ for $1 \leq i \leq N$ be units from one region where pooling operation is applied. Using Definition 2.2, for all $x \geq 0$ and some constant $K > 0$ we have

$$\mathbb{P}(|X_i| \geq x) \leq \exp\left(-x^{1/\theta}/K\right) \text{ for all } i.$$

Max pooling operation takes the maximum element in the region. Since X_i , $1 \leq i \leq N$ are the elements in one

region, we want to check if the tail of $\max_{1 \leq i \leq N} X_i$ obeys sub-Weibull property with optimal tail parameter equals to θ . Since max pooling operation can be decomposed into linear and ReLU operations, which does not harm the distribution tail (the extended envelope property for ReLU and Proposition 1.4 from Supplementary material), it leads to the proposition statement first part.

Summation and division by a constant does not influence the distribution tail (Proposition 1.4 from Supplementary material), yielding the proposition result regarding the averaging operation. \square

Corollary 2.1. *Consider a convolutional neural network containing convolutional, pooling and fully-connected layers under assumptions from Section 2.1. Then a unit of ℓ -th hidden layer has sub-Weibull distribution with optimal tail parameter $\theta = \ell/2$, where ℓ is the number of convolutional and fully-connected layers.*

Proof. Proposition 2.2 implies that the pooling layer keeps the tail parameter. From discussion at the beginning of the section, the result of Theorem 2.1 is also applied to convolutional neural networks where the depth is considered as the number of convolutional and fully-connected layers. \square

3. Regularization scheme on the units

Our main theoretical contribution, Theorem 2.1, characterizes the marginal prior distribution of the network units as follows: when the depth increases, the distribution becomes more heavy-tailed. In this section, we provide an interpretation of the result in terms of regularization at the level of the units. To this end, we first briefly recall shrinkage and penalized estimation methods.

3.1. Short digest on penalized estimation

The notion of penalized estimation is probably best illustrated on the simple linear regression model, where the aim is to improve prediction accuracy by shrinking, or even putting exactly to zero, some coefficients in the regression. Under these circumstances, inference is also more *interpretable* since, by reducing the number of coefficients effectively used in the model, it is possible to grasp its salient features. Shrinking is performed by imposing a penalty on the size of the coefficients, which is equivalent to allowing for a given budget on their size. Denote the regression parameter by $\beta \in \mathbb{R}^p$, the regression sum-of-squares by $R(\beta)$, and the penalty by $\lambda L(\beta)$, where L is some norm on \mathbb{R}^p and λ some positive tuning parameter. Then, the two

formulations of the regularized problem

$$\begin{aligned} & \min_{\beta \in \mathbb{R}^p} R(\beta) + \lambda L(\beta), \text{ and} \\ & \min_{\beta \in \mathbb{R}^p} R(\beta) \text{ subject to } L(\beta) \leq t, \end{aligned}$$

are equivalent, with some one-to-one correspondence between λ and t , and are respectively termed the *penalty* and the *constraint* formulation. This latter formulation provides an interesting geometrical intuition of the shrinkage mechanism: the constraint $L(\beta) \leq t$ reads as imposing a total budget of t for the parameter size in terms of the norm L . If the ordinary least squares estimator $\hat{\beta}^{\text{ols}}$ lives in the L -ball with surface $L(\beta) = t$, then there is no effect on the estimation. In contrast, when $\hat{\beta}^{\text{ols}}$ is outside the ball, then the intersection of the lowest level curve of the sum-of-squares $R(\beta)$ with the L -ball defines the penalized estimator.

The choice of the L norm has considerable effects on the problem, as can be sensed geometrically. Consider for instance \mathcal{L}^q norms, with $q \geq 0$. For any $q > 1$, the associated \mathcal{L}^q norm is differentiable and contours have a round shape without sharp angles. In that case, the penalty effect is to shrink the β coefficients towards 0. The most well-known estimator falling in this class is the *ridge* regression obtained with $q = 2$, see Figure 2 top-left panel. In contrast, for any $q \in (0, 1]$, the \mathcal{L}^q norm has some non differentiable points along the axis coordinates, see Figure 2 top-right and bottom panels. Such critical points are more likely to be hit by the level curves of the sum-of-squares $R(\beta)$, thus setting exactly to zero some of the parameters. A very successful approach in this class is the Lasso obtained with $q = 1$. Note that the problem is computationally much easier in the convex situation which occurs only for $q \geq 1$.

3.2. MAP on weights W is weight decay

These penalized methods have a simple Bayesian counterpart in the form of the maximum a posteriori (MAP) estimator. In this context, the objective function R is the negative log-likelihood, while the penalty L is the negative log-prior. The objective function takes on the form of sum-of-squared errors for regression under Gaussian errors, and of cross-entropy for classification.

For neural networks, it is well-known that an independent Gaussian prior on the weights

$$\pi(\mathbf{W}) \propto \prod_{\ell=1}^L \prod_{i,j} e^{-\frac{1}{2}(W_{i,j}^{(\ell)})^2},$$

is equivalent to the weight decay penalty, also known as ridge regression:

$$L(\mathbf{W}) = \sum_{\ell=1}^L \sum_{i,j} (W_{i,j}^{(\ell)})^2 = \|\mathbf{W}\|_2^2,$$

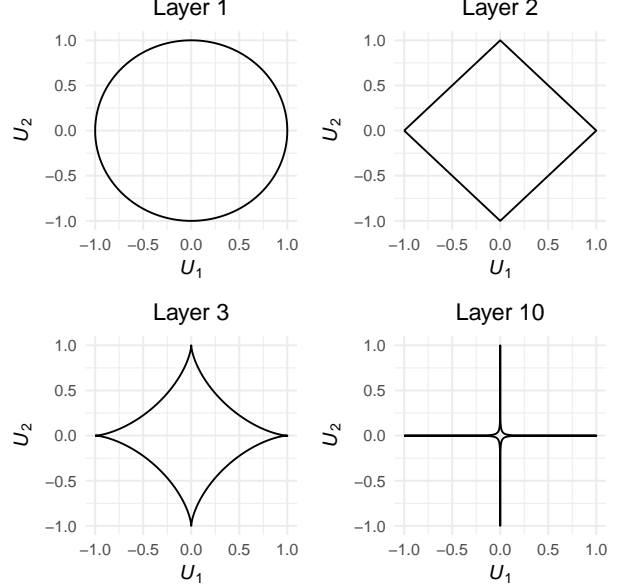


Figure 2. $\mathcal{L}^{2/\ell}$ -norm unit balls (in dimension 2) for layers $\ell = 1, 2, 3$ and 10.

where products and sums involving i and j above are over $1 \leq i \leq H_{\ell-1}$ and $1 \leq j \leq H_{\ell}$, H_0 and H_L representing respectively the input and output dimensions.

3.3. MAP on units U

Now moving the point of view from *weights* to *units* leads to a radically different shrinkage effect. Let $U_m^{(\ell)}$ denote the m -th unit of the ℓ -th layer (either pre- or post-nonlinearity). We prove in Theorem 2.1 that conditional on the input \mathbf{x} , a Gaussian prior on the weights translates into some prior on the units $U_m^{(\ell)}$ that is marginally sub-Weibull with optimal tail index $\theta = \ell/2$. This means that the tails of $U_m^{(\ell)}$ satisfy

$$\mathbb{P}(|U_m^{(\ell)}| \geq u) \leq \exp\left(-u^{2/\ell}/K_1\right) \quad \text{for all } u \geq 0, \quad (11)$$

for some positive constant K_1 . The exponent of u in the exponential term above is optimal in the sense that Equation (11) is not satisfied with some parameter θ' smaller than $\ell/2$. Thus, the marginal density of $U_m^{(\ell)}$ on \mathbb{R} is approximately proportional to

$$\pi_m^{(\ell)}(u) \approx e^{-|u|^{2/\ell}/K_1}.$$

The joint prior distribution for all the units $\mathbf{U} = (U_m^{(\ell)})_{1 \leq \ell \leq L, 1 \leq m \leq H_{\ell}}$ can be expressed from all the marginal distributions by Sklar's representation theorem (Sklar, 1959) as

$$\pi(\mathbf{U}) = \prod_{\ell=1}^L \prod_{m=1}^{H_{\ell}} \pi_m^{(\ell)}(U_m^{(\ell)}) C(F(\mathbf{U})), \quad (12)$$

where C represents the copula of \mathbf{U} (which characterizes all the dependence between the units) while F denotes its cumulative distribution function. The penalty incurred by such a prior distribution is obtained as the negative log-prior,

$$\begin{aligned} L(\mathbf{U}) &= -\sum_{\ell=1}^L \sum_{m=1}^{H_\ell} \log \pi_m^{(\ell)}(U_m^{(\ell)}) - \log C(F(\mathbf{U})), \\ &\approx \sum_{\ell=1}^L \sum_{m=1}^{H_\ell} |U_m^{(\ell)}|^{2/\ell} - \log C(F(\mathbf{U})), \\ &\approx \|\mathbf{U}^{(1)}\|_2^2 + \|\mathbf{U}_1^{(2)}\|_1 + \dots + \|\mathbf{U}^{(L)}\|_{2/L}^{2/L} \\ &\quad - \log C(F(\mathbf{U})). \end{aligned} \quad (13)$$

The first L terms in (13) indicate that some shrinkage operates at every layer of the network, with a penalty term that approximately takes the form of the $\mathcal{L}^{2/\ell}$ norm. Thus, the deeper the layer, the stronger the regularization induced at the level of the units, as summarized in Table 1.

Layer	Penalty on \mathbf{W}	Penalty on \mathbf{U}
1	$\ \mathbf{W}^{(1)}\ _2^2, \mathcal{L}^2$	$\ \mathbf{U}^{(1)}\ _2^2, \mathcal{L}^2$ (weight decay)
2	$\ \mathbf{W}^{(2)}\ _2^2, \mathcal{L}^2$	$\ \mathbf{U}^{(2)}\ _1, \mathcal{L}^1$ (Lasso)
ℓ	$\ \mathbf{W}^{(\ell)}\ _2^2, \mathcal{L}^2$	$\ \mathbf{U}^{(\ell)}\ _{2/\ell}^{2/\ell}, \mathcal{L}^{2/\ell}$

Table 1. Comparison of Bayesian neural network shrinkage effect on weights \mathbf{W} and units \mathbf{U} .

4. Experiments

We illustrate the result of Theorem 2.1 on a 100 layers MLP. The hidden layers of neural network have $H_1 = 1000$, $H_2 = 990$, $H_3 = 980$, \dots , $H_\ell = 1010 - 10\ell$, \dots , $H_{100} = 10$ hidden units, respectively. The input \mathbf{x} is a vector of features from \mathbb{R}^{10^4} . Figure 3 represents the tails of first three, 10th and 100th hidden layers pre-nonlinearity marginal distributions in logarithmic scale. The curves are obtained as histograms from a sample of size 10^5 from the prior on the pre-nonlinearitys, which is itself obtained by sampling 10^5 sets of weights \mathbf{W} from the Gaussian prior (2) and forward propagation via (1). The input vector \mathbf{x} is sampled with independent features from a standard normal distribution once for all at the start. The nonlinearity ϕ is the ReLU function. Being a linear combination involving symmetric weights \mathbf{W} , pre-nonlinearitys \mathbf{g} also have a symmetric distribution, thus we visualize only their distribution on \mathbb{R}_+ .

Figure 3 corroborates our main result. On the one hand, the prior distribution of the first hidden units is Gaussian (green curve), which corresponds to a **subW**(1/2) distribution. On

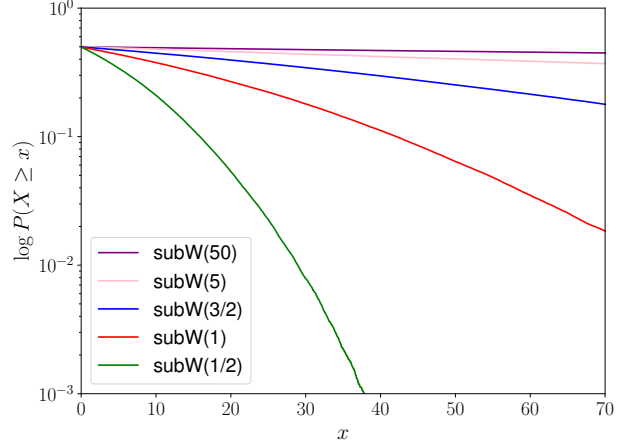


Figure 3. Illustration of layers $\ell = 1, 2, 3, 10$ and 100 hidden units (pre-nonlinearitys) marginal prior distributions. They correspond respectively to **subW**(1/2), **subW**(1), **subW**(3/2), **subW**(5) and **subW**(50).

the other hand, deeper layers are characterized by heavier-tailed distributions. The deepest considered layer (100th, violet curve) has an extremely flat distribution, which corresponds to a **subW**(50) distribution.

5. Conclusion and future work

Despite the ubiquity of deep learning throughout science, medicine and engineering, the underlying theory has not kept pace with applications for deep learning. In this paper, we have extended the state of knowledge on Bayesian neural networks by providing a characterization of the marginal prior distribution of the units. We proved that they are heavier-tailed as depth increases, and discussed this result in terms of a regularizing mechanism at the level of the units.

Since initialization and learning dynamics are key in modern machine learning in order to properly tune deep learning algorithms, a good implementation practice requires a proper understanding of the prior distribution at play and of the regularization it incurs.

We hope that our results will open avenues for further research. Firstly, Theorem 2.1 regards the *marginal* prior distribution of the units, while a full characterization of the joint distribution of all units \mathbf{U} remains an open question. More specifically, a precise description of the copula defined in Equation (12) would provide valuable information about the dependence between the units, and also about the precise geometrical structure of the balls induced by that penalty. Secondly, the interpretation of our result (Section 3) is concerned with the maximum a posteriori of the units, which is a point estimator. One of the benefits of the Bayesian approach to neural networks lies in its ability to

provide a principled approach to uncertainty quantification, so that an interpretation of our result in terms of the full posterior distribution would be very appealing. Lastly, the practical potentialities of our results are many: to better comprehend the regularizing mechanisms in deep neural networks will contribute to design and understand strategies to avoid overfitting and improve generalization.

Acknowledgements

We would like to thank [Stéphane Girard](#) for fruitful discussions on Weibull-like distributions and [Cédric Févotte](#) for pointing out the potential relationship of our heavy-tail result with sparsity-inducing priors. This work has been partially supported by the grant ANR-16-CE23-0006 Deep in France and LabEx PERSYVAL-Lab (ANR-11-LABX-0025-01).

References

- Ba, J. L., Kiros, J. R., and Hinton, G. E. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- Bibi, A., Alfadly, M., and Ghanem, B. Analytic Expressions for Probabilistic Moments of PL-DNN with Gaussian Input. In *The IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- Cho, Y. and Saul, L. K. Kernel methods for deep learning. In *Advances in Neural Information Processing Systems*, pp. 342–350, 2009.
- Damianou, A. and Lawrence, N. Deep Gaussian processes. In *Artificial Intelligence and Statistics*, pp. 207–215, 2013.
- Duvenaud, D., Rippel, O., Adams, R., and Ghahramani, Z. Avoiding pathologies in very deep networks. In *Artificial Intelligence and Statistics*, pp. 202–210, 2014.
- Fukushima, K. and Miyake, S. Neocognitron: A self-organizing neural network model for a mechanism of visual pattern recognition. In *Competition and cooperation in neural nets*, pp. 267–285. Springer, 1982.
- Gal, Y. and Ghahramani, Z. Dropout as a Bayesian approximation: Representing model uncertainty in deep learning. In *International Conference on Machine Learning*, pp. 1050–1059, 2016.
- Goodfellow, I., Bengio, Y., and Courville, A. *Deep learning*. MIT press, 2016.
- Graves, A., Mohamed, A., and Hinton, G. E. Speech recognition with deep recurrent neural networks. In *International Conference on Acoustics, Speech and Signal Processing*, pp. 6645–6649, 2013.
- Hron, J., Matthews, A. G. d. G., and Ghahramani, Z. Variational Bayesian dropout: pitfalls and fixes. *arXiv preprint arXiv:1807.01969*, 2018.
- Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In Bach, F. and Blei, D. (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 448–456, Lille, France, 07–09 Jul 2015. PMLR.
- Kingma, D. P., Salimans, T., and Welling, M. Variational dropout and the local reparameterization trick. In *Advances in Neural Information Processing Systems*, pp. 2575–2583, 2015.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. In *Neural Information Processing Systems*, pp. 1097–1105, 2012.
- Krogh, A. and Hertz, J. A. A simple weight decay can improve generalization. In *Neural Information Processing Systems*, pp. 950–957, 1991.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Lee, J., Sohl-dickstein, J., Pennington, J., Novak, R., Schoenholz, S., and Bahri, Y. Deep neural networks as Gaussian processes. In *International Conference on Learning Representations*, 2018.
- Liang, F., Li, Q., and Zhou, L. Bayesian neural networks for selection of drug sensitive genes. *Journal of the American Statistical Association*, 113(523):955–972, 2018.
- MacKay, D. J. A practical Bayesian framework for backpropagation networks. *Neural Computation*, 4(3):448–472, 1992.
- Matthews, A. G. d. G., Rowland, M., Hron, J., Turner, R. E., and Ghahramani, Z. Gaussian process behaviour in wide deep neural networks. *arXiv preprint arXiv:1804.11271*, 2018a.
- Matthews, A. G. d. G., Rowland, M., Hron, J., Turner, R. E., and Ghahramani, Z. Gaussian process behaviour in wide deep neural networks. In *Proceedings of the 6th International Conference on Learning Representations*, 2018b.
- Molchanov, D., Ashukha, A., and Vetrov, D. Variational dropout sparsifies deep neural networks. In *International Conference on Machine Learning*, pp. 2498–2507, 2017.
- Neal, R. M. Bayesian training of backpropagation networks by the hybrid Monte Carlo method. Technical report, Citeseer, 1992.

- Neal, R. M. *Bayesian learning for neural networks*, volume 118. Springer Science & Business Media, 1996.
- Polson, N. G. and Sokolov, V. Deep learning: A Bayesian perspective. *Bayesian Analysis*, 12(4):1275–1304, 2017.
- Rasmussen, C. E. and Williams, C. K. I. *Gaussian Processes for Machine Learning*. MIT Press, 2006. doi: 10.1.1.86.3414.
- Rinne, H. *The Weibull distribution: a handbook*. Chapman and Hall/CRC, 2008.
- Saatci, Y. and Wilson, A. G. Bayesian GAN. In *Advances in Neural Information Processing Systems*, pp. 3622–3631, 2017.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., and Lanctot, M. Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- Sklar, M. Fonctions de repartition an dimensions et leurs marges. *Publ. inst. statist. univ. Paris*, 8:229–231, 1959.
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- Tibshirani, R. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pp. 267–288, 1996.

A. Equivalent sub-Weibull distribution properties

Proposition A.1 (Sub-Weibull distribution). *Let X be a random variable. Then the following properties are equivalent; the parameters $K_i > 0$ appearing in these properties differ from each other by at most an absolute constant factor.*

1. The tails of X satisfy

$$\mathbb{P}(|X| \geq x) \leq 2 \exp\left(-x^{1/\theta}/K_1\right) \quad \text{for all } x \geq 0.$$

2. The moments of X satisfy

$$\|X\|_k = (\mathbb{E}[|X|^k])^{1/k} \leq K_2 k^\theta \quad \text{for all } k \geq 1.$$

3. The MGF of $X^{1/\theta}$ satisfies

$$\mathbb{E}\left[\exp\left(\lambda^{1/\theta} X^{1/\theta}\right)\right] \leq K_2 \exp(K_3^{1/\theta} \lambda^{1/\theta})$$

for all λ such that $|\lambda| \leq \frac{1}{K_3}$.

4. The MGF of $X^{1/\theta}$ is bounded at some point, namely

$$\mathbb{E}\left[\exp\left(X^{1/\theta}/K_4\right)\right] \leq 2.$$

Proof. **1** \Rightarrow **2**. Assume property **1** holds. Applying the integral identity for $|X|^k$, we obtain

$$\begin{aligned} \mathbb{E}[|X|^k] &= \int_0^\infty \mathbb{P}(|X|^k > x) dx \\ &= \int_0^\infty \mathbb{P}(|X| > x^{1/k}) dx \\ &\leq \int_0^\infty 2 \exp\left(-x^{1/(k\theta)}/K_1\right) dx \\ &= 2K_1^{k\theta} k\theta \int_0^\infty e^{-u} u^{k\theta-1} du = 2K_1^{k\theta} k\theta \Gamma(k\theta) \\ &\sim K_1^{k\theta} k\theta (k\theta - 1)^{k\theta-1} \sim (K_1 k\theta)^{k\theta}. \end{aligned}$$

Taking the k -th root of the expression above yields property **2**

$$\|X\|_k \lesssim (K_1 \theta)^\theta k^\theta \leq K_2 k^\theta,$$

with $K_2 = (K_1 \theta)^\theta$.

2 \Rightarrow **3**. Assume property **2** holds. Recalling the Taylor series expansion of the exponential function, we obtain

$$\begin{aligned} \mathbb{E}\left[\exp\left(\lambda^{1/\theta} X^{1/\theta}\right)\right] &= \mathbb{E}\left[1 + \sum_{k=1}^\infty \frac{(\lambda^{1/\theta} |X|^{1/\theta})^k}{k!}\right] \\ &= 1 + \sum_{k=1}^\infty \frac{\lambda^{k/\theta} \mathbb{E}[|X|^{k/\theta}]}{k!}. \end{aligned}$$

Property **2** guarantees that $\mathbb{E}[|X|^k] \leq K_2 k^{k/\theta}$ and $\mathbb{E}[|X|^{k/\theta}] \leq K_2 (k/\theta)^k$ for some K_2 . Stirling's approximation yields $k! \geq (k/e)^k$. Substituting these two bounds, we get

$$\begin{aligned} \mathbb{E}\left[\exp\left(\lambda^{1/\theta} X^{1/\theta}\right)\right] &\leq \sum_{k=1}^\infty \frac{\lambda^{k/\theta} K_2 (k/\theta)^k}{(k/e)^k} \\ &= \sum_{k=0}^\infty K_2 (e\lambda^{1/\theta}/\theta)^k = \frac{K_2}{1 - e\lambda^{1/\theta}/\theta}, \end{aligned}$$

provided that $e\lambda^{1/\theta}/\theta < 1$, in which case the geometric series above converges. To bound this quantity further, we can use the numeric inequality $\frac{1}{1-x} \leq e^{2x}$, which is valid for $x \in [0, 1/2]$. It follows that

$$\mathbb{E}\left[\exp\left(\lambda^{1/\theta} X^{1/\theta}\right)\right] \leq K_2 \exp\left(2e\lambda^{1/\theta}/\theta\right)$$

for all λ satisfying $|\lambda| \leq (\frac{\theta}{2e})^\theta$. This yields property **3** with $K_3 = (2e/\theta)^\theta$.

3 \Rightarrow **4**. Assume property **3** holds. Take $\lambda = 1/K_4$, where $K_4 \geq K_3/(\ln 2 - \ln K_2)^\theta$. This yields property **4**.

4 \Rightarrow **1**. Assume property **4** holds. We may assume that $K_4 = 1$. Then, by Markov's inequality and property **3**, we obtain

$$\begin{aligned} \mathbb{P}(|X| > x) &= \mathbb{P}(e^{|X|^{1/\theta}} > e^{x^{1/\theta}}) \\ &\leq \frac{\mathbb{E}[e^{|X|^{1/\theta}}]}{e^{x^{1/\theta}}} \leq 2e^{-x^{1/\theta}/K_1}. \end{aligned}$$

This proves property **1** with $K_1 = 1$. \square

Remark A.1. The constant 2 that appears in some properties in Proposition A.1 does not have any special meaning. It is chosen for simplicity and can be replaced by other absolute constants.

Definition A.1 (Sub-Weibull random variable). *A random variable X that satisfies one of the equivalent properties of Proposition A.1 is called a sub-Weibull random variable with tail parameter θ , which is denoted by $X \sim \text{subW}(\theta)$.*

Informally, the tails of a $\text{subW}(\theta)$ distribution are dominated by (i.e. decay at least as fast as) the tails of a Weibull variable with the shape parameter equal to $1/\theta$ (Rinne, 2008). The larger tail parameter θ , the heavier the tails of the sub-Weibull distribution.

Sub-Gaussian and sub-Exponential variables, which are commonly used, are special cases of sub-Weibull random variables with tail parameter $\theta = 1/2$ and $\theta = 1$, respectively (see Table 2).

Distribution	Tail	Moments
Sub-Gaussian	$\bar{F}(x) \leq e^{-\lambda x^2}$	$\ X\ _k \leq C\sqrt{k}$
Sub-Exponential	$\bar{F}(x) \leq e^{-\lambda x}$	$\ X\ _k \leq Ck$
Sub-Weibull	$\bar{F}(x) \leq e^{-\lambda x^{1/\theta}}$	$\ X\ _k \leq Ck^\theta$

Table 2. Sub-Gaussian, sub-Exponential and sub-Weibull distributions comparison in terms of tail $\bar{F}(x) = P(X \geq x)$ and moment condition, with λ and C some positive constants. The first two are a special case of the last with $\theta = 1/2$ and $\theta = 1$ respectively.

Proposition A.2 (Inclusion). *Let θ_1 and θ_2 such that $0 < \theta_1 < \theta_2$ be tail parameters for some sub-Weibull distributed variables. Then the following inclusion holds*

$$\text{subW}(\theta_1) \subset \text{subW}(\theta_2).$$

Proof. For $X \sim \text{subW}(\theta_1)$, it holds that $\|X\|_k \leq K_2 k^{\theta_1}$. Since $k^{\theta_1} \leq k^{\theta_2}$ for all $k \geq 1$, this yields $\|X\|_k \leq K_2 k^{\theta_2}$, which by definition implies $X \sim \text{subW}(\theta_2)$. \square

The following proposition is key in establishing that neural network units of layer ℓ are $\text{subW}(\ell/2)$, where $\ell/2$ is optimal.

Proposition A.3 (Optimal sub-Weibull tail coefficient and moment condition). *Let $\theta > 0$ and let X be a random variable satisfying the following asymptotic equivalence on moments*

$$\|X\|_k \asymp k^\theta.$$

Then X is sub-Weibull distributed with optimal tail parameter θ , in the sense that for any $\theta' < \theta$, X is not sub-Weibull with tail parameter θ' .

Proof. Since X satisfies Condition 2 of Proposition A.1, $X \sim \text{subW}(\theta)$. Let $\theta' < \theta$. Since $\|X\|_k \asymp k^\theta$, there does not exist any constant K_2 such that $\|X\|_k \leq K_2 k^{\theta'}$, so X is not sub-Weibull with tail parameter θ' . \square

Proposition A.4 (Elementary operations). *Summation: If X and Y are independent and from the one sub-Weibull distribution family, then $X + Y$ belongs to the same sub-Weibull distribution family: $X, Y \sim \text{subW}(\theta)$ implies $X + Y \sim \text{subW}(\theta)$.*

Multiplication by a constant: X and cX are from the one sub-Weibull distribution family, where $c > 0$ is a constant.

Proof. Let $X, Y \sim \text{subW}(\theta)$, then for all $x \geq 0$ and some constant $K > 0$ we have from Proposition A.1, property 3:

$$\begin{aligned} \mathbb{E} \left[\exp \left(\lambda^{1/\theta} (X + Y)^{1/\theta} \right) \right] &= \\ &= \mathbb{E} \left[\exp \left(\lambda^{1/\theta} X^{1/\theta} \right) \right] \mathbb{E} \left[\exp \left(\lambda^{1/\theta} Y^{1/\theta} \right) \right] \leq \\ &\leq K_{2X} \exp(K_{3X}^{1/\theta} \lambda^{1/\theta}) \cdot \leq K_{2Y} \exp(K_{3Y}^{1/\theta} \lambda^{1/\theta}) = \\ &= K_2 \exp(K_3^{1/\theta} \lambda^{1/\theta}). \end{aligned}$$

And multiplication by a constant we obtain from Proposition A.1, property 1:

$$\mathbb{P}(CX \geq x) = \mathbb{P}(X \geq x/C) \leq \exp \left(-x^{1/\theta}/K \right).$$

\square

It is typically assumed that the random variable X has zero mean. If this is not the case, we can always center X by subtracting the mean, not changing the tail parameter of sub-Weibull distribution it follows.

Lemma A.1 (Centered variables). *Centering does not harm tail properties. In particular, random variables X and $(X - \mathbb{E}[X])$ belong to the same sub-Weibull family, i.e. with the same optimal tail parameter.*

Proof. Let prove that $\|X\|_k \asymp \|X - \mathbb{E}[X]\|_k$. Consider $\|X - \mathbb{E}[X]\|_k$. According to triangle inequality, we have

$$\|X - \mathbb{E}[X]\|_k \leq \|X\|_k + \|\mathbb{E}[X]\|_k.$$

Since $\mathbb{E}[X]$ equals to some constant, then $\|\mathbb{E}[X]\|_k = \|\mathbb{E}[X]\|_1 \leq \mathbb{E}[|X|] = \|X\|_1$.

Let show that the norm is increasing function with respect to k . For real-valued random variables Y and Z Holder's inequality reads

$$\mathbb{E}[|YZ|] \leq (\mathbb{E}[|Y|^p])^{1/p} (\mathbb{E}[|Z|^q])^{1/q}.$$

Let $0 < r < s$ and define $p = s/r$. Then $q = p/(p-1)$ is the Holder conjugate of p . Applying Holder's inequality to the random variable $|Y|^r$, obtain

$$\mathbb{E}[|Y|^r] \leq (\mathbb{E}[|Y|^s])^{r/s}.$$

Taking r -th roots results in the inequality $\|Y\|_r \leq \|Y\|_s$ for $r < s$.

Hence $\|X\|_1 \leq C_0 \|X\|_k$ with $C_0 > 0$ and

$$\|X - \mathbb{E}[X]\|_k \leq (C_0 + 1) \|X\|_k.$$

Consider $\|X\|_k$. According to triangle inequality, we have

$$\|X\|_k = \|X - \mathbb{E}[X] + \mathbb{E}[X]\|_k \leq \|X - \mathbb{E}[X]\|_k + \|\mathbb{E}[X]\|_k.$$

Since $\|\mathbb{E}[X]\|_k = \|X\|_1 \leq C_0 \|X\|_k$ for $C_0 > 0$ and the inequality $\|X\|_1 \leq \|X\|_k$ holds, choose the constant $C_0 < 1$ such that $\|X\|_1 \leq C_0 \|X\|_k \leq \|X\|_k$. Then we have

$$\|X\|_k \leq \|X - \mathbb{E}[X]\|_k + C_0 \|X\|_k.$$

As $1 - C_0 > 0$, we obtain

$$(1 - C_0) \|X\|_k \leq \|X - \mathbb{E}[X]\|_k.$$

It implies $\|X\|_k \asymp \|X - \mathbb{E}[X]\|_k$.

Let $X \sim \text{subW}(\theta)$, then $\|X\|_k \leq Ck^\theta$. Due to asymptotic equivalence, the inequality $\|X - \mathbb{E}[X]\| \leq \tilde{C}k^\theta$ holds. This is sufficient condition for $(X - \mathbb{E}[X]) \sim \text{subW}(\theta)$. Analogously, by assuming belonging $(X - \mathbb{E}[X])$ to sub-Weibull with θ , obtain that the variable X follows sub-Weibull with θ . This ends the proof. \square

B. Intermediate lemmas

Introduce the definition of asymptotic equivalence between numeric sequences:

Definition B.1 (Asymptotic equivalence). *Two sequences a_k and b_k are called asymptotic equivalent and denoted as $a_k \asymp b_k$ if there exist constants $d > 0$ and $D > 0$ such that*

$$d \leq \frac{a_k}{b_k} \leq D, \quad \text{for all } k \in \mathbb{N}. \quad (14)$$

Lemma B.1 (Gaussian moments). *Let X be a normal random variable such that $X \sim \mathcal{N}(0, \sigma^2)$, then the following asymptotic equivalence holds*

$$\|X\|_k \asymp \sqrt{k}.$$

Proof. The moments of central normal absolute random variable $|X|$ are equal to

$$\begin{aligned} \mathbb{E}[|X|^k] &= \int_{\mathbb{R}} |x|^k p(x) dx \\ &= 2 \int_0^\infty x^k p(x) dx \\ &= \frac{1}{\sqrt{\pi}} \sigma^k 2^{k/2} \Gamma\left(\frac{k+1}{2}\right). \end{aligned} \quad (15)$$

We have the expression for the Gamma function

$$\Gamma(z) = \sqrt{\frac{2\pi}{z}} \left(\frac{z}{e}\right)^z \left(1 + \frac{1}{12z} + o\left(\frac{1}{z}\right)\right). \quad (16)$$

Substituting (16) into the central normal absolute moment (15), we obtain

$$\begin{aligned} \mathbb{E}[|X|^k] &= \frac{1}{\sqrt{\pi}} \sigma^k 2^{k/2} \sqrt{\frac{4\pi}{k+1}} \left(\frac{k+1}{2e}\right)^{(k+1)/2} \\ &\quad \cdot \left(1 + \frac{1}{6(k+1)} + o\left(\frac{1}{k}\right)\right) \\ &= \frac{2\sigma^k}{\sqrt{2e}} \left(\frac{k+1}{e}\right)^{k/2} \left(1 + \frac{1}{6(k+1)} + o\left(\frac{1}{k}\right)\right). \end{aligned}$$

Then the roots of absolute moments can be written in the form of

$$\begin{aligned} \|X\|_k &= \frac{\sigma}{e^{1/(2k)}} \sqrt{\frac{k+1}{e}} \left(1 + \frac{1}{6(k+1)} + o\left(\frac{1}{k}\right)\right)^{1/k} \\ &= \frac{\sigma}{e} \frac{\sqrt{k+1}}{e^{1/(2k)}} \left(1 + \frac{1}{6(k+1)k} + o\left(\frac{1}{k^2}\right)\right) \\ &= \frac{\sigma}{e} c_k \sqrt{k+1}. \end{aligned}$$

Here the coefficient c_k denotes

$$c_k = \frac{1}{e^{1/(2k)}} \left(1 + \frac{1}{6(k+1)k} + o\left(\frac{1}{k^2}\right)\right) \rightarrow 1,$$

with $k \rightarrow \infty$. Thus, asymptotic equivalence holds

$$\|X\|_k \asymp \sqrt{k+1} \asymp \sqrt{k}.$$

\square

Lemma B.2 (Multiplication moments). *Let W and X be independent random variables such that $W \sim \mathcal{N}(0, \sigma^2)$ and for some $p > 0$ it holds*

$$\|X\|_k \asymp k^p. \quad (17)$$

Let W_i be independent copies of W , and X_i be copies of X , $i = 1, \dots, H$ with non-negative covariance between moments of copies

$$\text{Cov}[X_i^s, X_j^t] \geq 0, \quad \text{for } i \neq j, s, t \in \mathbb{N}. \quad (18)$$

Then we have the following asymptotic equivalence

$$\left\| \sum_{i=1}^H W_i X_i \right\|_k \asymp k^{p+1/2}. \quad (19)$$

Proof. Let us proof the statement, using mathematical induction.

Base case: show that the statement is true for $H = 1$. For independent variables W and X , we have

$$\begin{aligned} \|WX\|_k &= (\mathbb{E}[|WX|^k])^{1/k} = (\mathbb{E}[|W|^k] \mathbb{E}[|X|^k])^{1/k} \\ &= \|W\|_k \|X\|_k. \end{aligned} \quad (20)$$

Since the random variable W follows Gaussian distribution, then Lemma B.1 implies

$$\|W\|_k \asymp \sqrt{k}. \quad (21)$$

Substituting assumption (17) and weight norm asymptotic equivalence (21) into (20) leads to the desired asymptotic equivalence (19) in case of $H = 1$.

Inductive step: show that if for $H = n - 1$ the statement holds, then for $H = n$ it also holds.

Suppose for $H = n - 1$ we have

$$\left\| \sum_{i=1}^{n-1} W_i X_i \right\|_k \asymp k^{p+1/2}. \quad (22)$$

Then, according to the covariance assumption (18), for $H = n$ we get

$$\begin{aligned} \left\| \sum_{i=1}^n W_i X_i \right\|_k &= \left\| \sum_{i=1}^{n-1} W_i X_i + W_n X_n \right\|_k \\ &\geq \sum_{j=0}^k C_k^j \left\| \sum_{i=1}^{n-1} W_i X_i \right\|_j^j \left\| W_n X_n \right\|_{k-j}^{k-j}. \end{aligned} \quad (23)$$

$$(24)$$

Using the equivalence definition (Def. B.1), from the induction assumption (22) for all $j = 0, \dots, k$ there exists absolute constant $d_1 > 0$ such that

$$\left\| \sum_{i=1}^{n-1} W_i X_i \right\|_j^j \geq (d_1 j^{p+1/2})^j. \quad (25)$$

Recalling previous equivalence results in the base case, there exists constant $m_2 > 0$ such that

$$\left\| W_n X_n \right\|_{k-j}^{k-j} \geq (d_2 (k-j)^{p+1/2})^{k-j}. \quad (26)$$

Substitute obtained bounds (25) and (26) into equation (23) with denoted $d = \min\{d_1, d_2\}$, obtain

$$\begin{aligned} \left\| \sum_{i=1}^n W_i X_i \right\|_k^k &\geq d^k \sum_{j=0}^k C_k^j [j^j (k-j)^{k-j}]^{p+1/2} \\ &= d^k k^{k(p+1/2)} \sum_{j=0}^k C_k^j \left[\left(\frac{j}{k} \right)^j \left(1 - \frac{j}{k} \right)^{k-j} \right]^{p+1/2}. \end{aligned} \quad (27)$$

Notice the lower bound of the following expression

$$\begin{aligned} \sum_{j=0}^k C_k^j \left[\left(\frac{j}{k} \right)^j \left(1 - \frac{j}{k} \right)^{k-j} \right]^{p+1/2} \\ \geq \sum_{j=0}^k \left[\left(\frac{j}{k} \right)^j \left(1 - \frac{j}{k} \right)^{k-j} \right]^{p+1/2} \geq 2. \end{aligned} \quad (28)$$

Substituting found lower bound (28) into (27), get

$$\left\| \sum_{i=1}^n W_i X_i \right\|_k^k \geq 2 d^k k^{k(p+1/2)} > d^k k^{k(p+1/2)}. \quad (29)$$

Now prove the upper bound. For random variables Y and Z the Holder's inequality holds

$$\begin{aligned} \|YZ\|_1 &= \mathbb{E}[|YZ|] \leq (\mathbb{E}[|Y|^2] \mathbb{E}[|Z|^2])^{1/2} \\ &= \|YZ\|_2 \|Y\|_2 \|Z\|_2. \end{aligned}$$

Holder's inequality leads to the inequality for L^k norm

$$\|YX\|_k^k \leq \|Y\|_{2k}^k \|X\|_{2k}^k. \quad (30)$$

Obtain the upper bound of $\left\| \sum_{i=1}^n W_i X_i \right\|_k^k$ from the norm property (30) for the random variables $Y = \left(\sum_{i=1}^{n-1} W_i X_i \right)^{k-j}$ and $Z = (W_n X_n)^j$

$$\begin{aligned} \left\| \sum_{i=1}^n W_i X_i \right\|_k^k &= \left\| \sum_{i=1}^{n-1} W_i X_i + W_n X_n \right\|_k^k \\ &\leq \sum_{j=0}^k C_k^j \left\| \sum_{i=1}^{n-1} W_i X_i \right\|_{2j}^j \left\| W_n X_n \right\|_{2(k-j)}^{k-j}. \end{aligned} \quad (31)$$

$$(32)$$

From the induction assumption (22) for all $j = 0, \dots, k$ there exists absolute constant $D_1 > 0$ such that

$$\left\| \sum_{i=1}^{n-1} W_i X_i \right\|_{2j}^j \leq (D_1 (2j)^{p+1/2})^j. \quad (33)$$

Recalling previous equivalence results in the base case, there exists constant $D_2 > 0$ such that

$$\left\| W_n X_n \right\|_{2(k-j)}^{k-j} \leq (D_2 (2(k-j))^{p+1/2})^{k-j}. \quad (34)$$

Substitute obtained bounds (33) and (34) into equation (31) with denoted $D = \max\{D_1, D_2\}$, obtain

$$\left\| \sum_{i=1}^n W_i X_i \right\|_k^k \leq D^k \sum_{j=0}^k C_k^j \left[(2j)^j (2(k-j))^{k-j} \right]^{p+1/2}.$$

Find an upper bound for $\left[\left(1 - \frac{j}{k}\right)^{k-j} \left(\frac{j}{k}\right)^j \right]^{p+1/2}$. Since expressions $\left(1 - \frac{j}{k}\right)$ and $\left(\frac{j}{k}\right)$ are less than 1, then $\left[\left(1 - \frac{j}{k}\right)^{k-j} \left(\frac{j}{k}\right)^j \right]^{p+1/2} < 1$ holds for all natural numbers $p > 0$. For the sum of binomial coefficients it holds the inequality $\sum_{j=0}^k C_k^j < 2^k$. So the final upper bound is

$$\left\| \sum_{i=1}^n W_i X_i \right\|_k^k \leq 2^k D^k (2k)^{k(p+1/2)}. \quad (35)$$

Hence, taking the k -th root of (29) and (35), we have upper and lower bounds which imply the equivalence for $H = n$ and the truth of inductive step

$$d' k^{p+1/2} \leq \left\| \sum_{i=1}^n W_i X_i \right\|_k \leq D' k^{p+1/2},$$

where $d' = d$ and $D' = 2^{p+3/2} D$. Since both the base case and the inductive step have been performed, by mathematical induction the equivalence holds for all $H \in \mathbb{N}$

$$\left\| \sum_{i=1}^H W_i X_i \right\|_k \asymp k^{p+1/2}.$$

□