

# Bayesian neural networks priors at the level of units

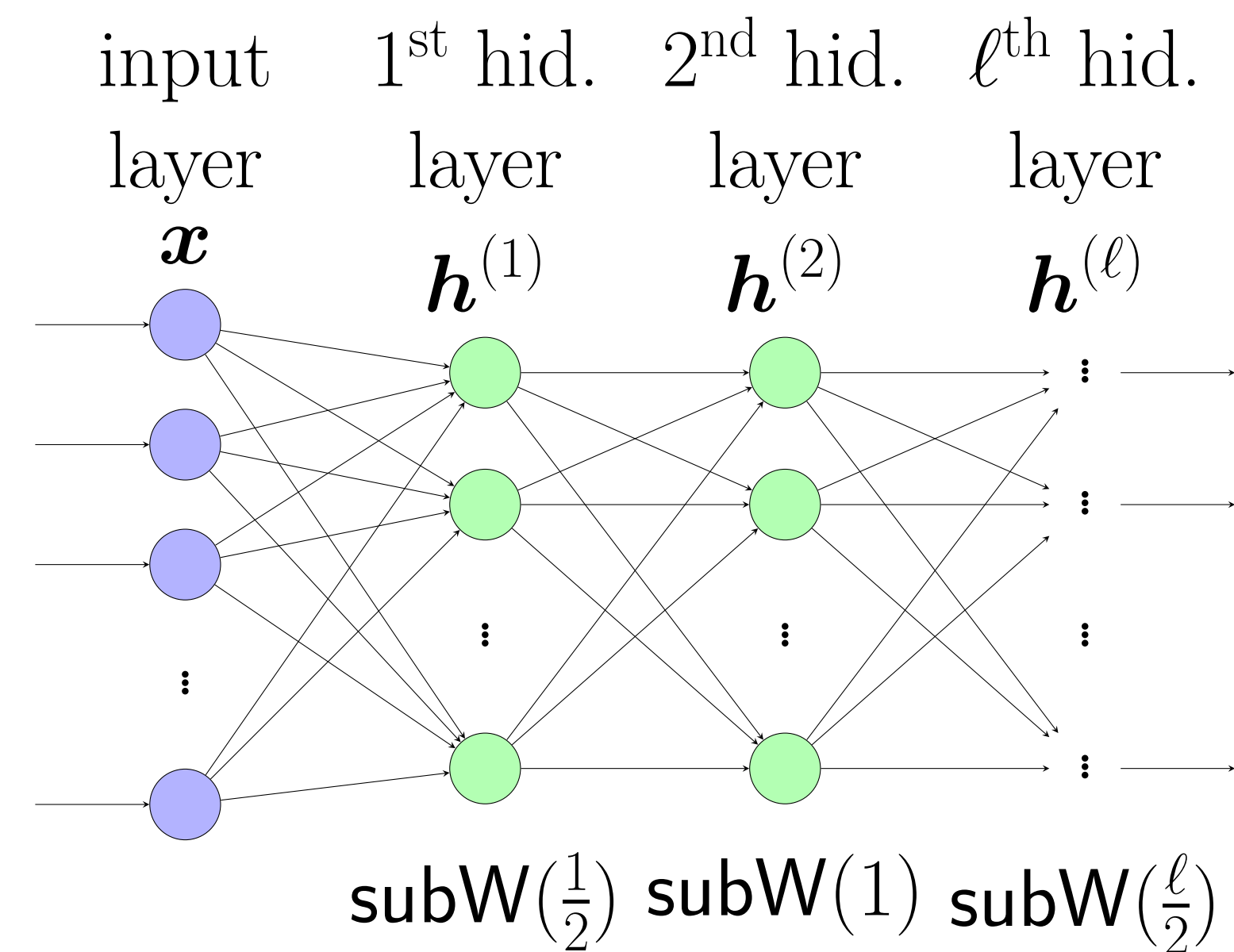
Mariia Vladimirova<sup>1</sup>, Julyan Arbel<sup>1</sup> and Pablo Mesejo<sup>2</sup>

<sup>1</sup>Inria Grenoble Rhône-Alpes, France

<sup>2</sup>University of Granada, Spain

## Introduction

We investigate deep Bayesian neural networks with Gaussian priors on the weights and ReLU-like nonlinearities. See Vladimirova et al. (2018).



## Notations

Given an input  $\mathbf{x} \in \mathbb{R}^N$ , the  $\ell$ -th hidden layer unit activations are defined as

$$\begin{aligned} \mathbf{g}^{(\ell)}(\mathbf{x}) &= \mathbf{W}^{(\ell)} \mathbf{h}^{(\ell-1)}(\mathbf{x}), \\ \mathbf{h}^{(\ell)}(\mathbf{x}) &= \phi(\mathbf{g}^{(\ell)}(\mathbf{x})). \end{aligned}$$

## Assumptions

- **Gaussian prior** on weights:

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

- A *nonlinearity*  $\phi: \mathbb{R} \rightarrow \mathbb{R}$  is said to obey the **extended envelope property** if there exist  $c_1, c_2, d_2 \geq 0, d_1 > 0$  such that

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

## Sub-Weibull

A random variable  $X$ , such that

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

for all  $x \geq 0$  and for some  $K > 0$ , is called a **sub-Weibull** random variable with tail parameter  $\theta > 0$ :

$$X \sim \text{subW}(\theta).$$

**Moment property:**

$X \sim \text{subW}(\theta)$  implies

$$\|X\|_k = (\mathbb{E}|X|^k)^{1/k} \asymp k^\theta,$$

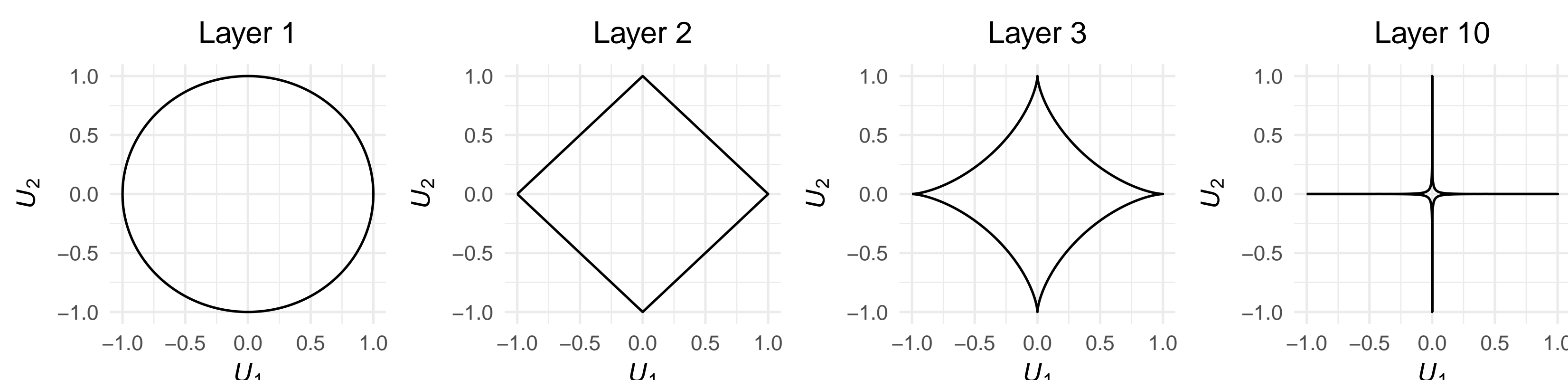
Meaning for all  $k \in \mathbb{N}$  and for some constants  $d, D > 0$ ,

$$d < \|X\|_k/k^\theta < D.$$

## Covariance theorem

The **covariance** between hidden units of the same layer is **non-negative**. Moreover, for any  $\ell$ -th hidden layer units  $h^{(\ell)}$  and  $\tilde{h}^{(\ell)}$ , for  $s, t \in \mathbb{N}$  it holds

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



## Theorem (Vladimirova et al., 2018)

The  $\ell$ -th hidden layer units  $U^{(\ell)}$  (pre-activation  $g^{(\ell)}$  or post-activation  $h^{(\ell)}$ ) of a feed-forward Bayesian neural network with:

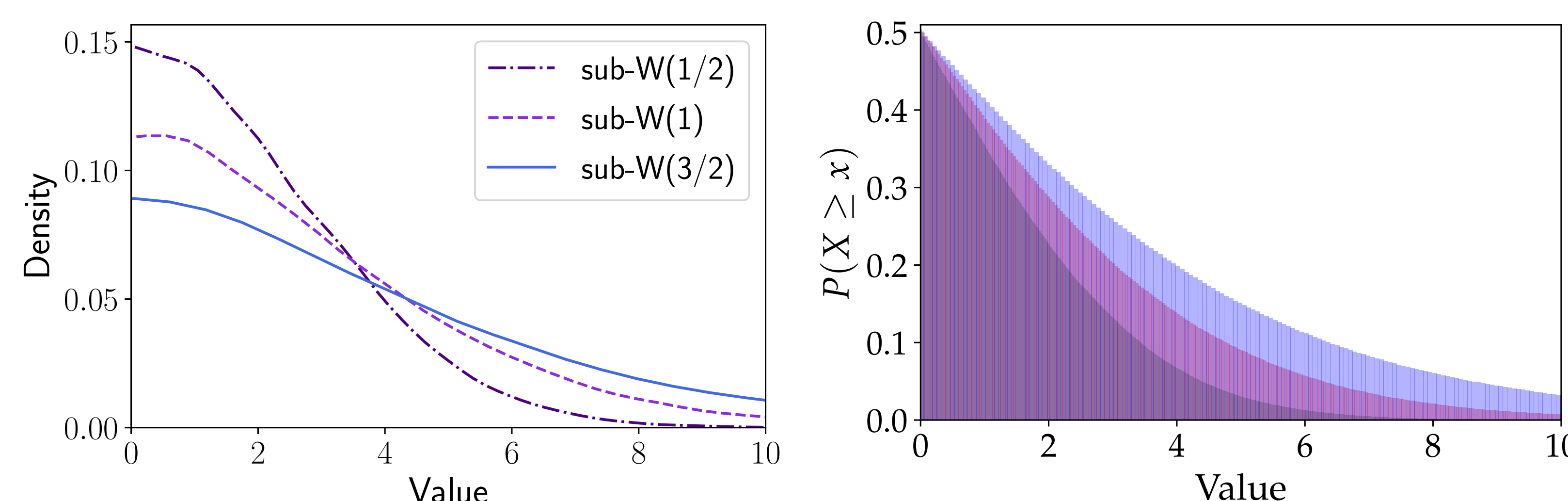
- **Gaussian priors** on weights and
- **extended envelope condition** activation function  $\phi$

have **sub-Weibull marginal prior distribution** with optimal tail parameter  $\theta = \ell/2$ , conditional on the input  $\mathbf{x}$ :

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

## Prior distributions of layers $\ell = 1, 2, 3$

Illustration of units marginal prior distributions from the first three hidden layers. Neural network parameters:  $(N, H_1, H_2, H_3) = (50, 25, 24, 4)$ .



## Proof sketch

Induction w.r.t. layer depth  $\ell$ :

$$\|h^{(\ell)}\|_k \asymp k^{\ell/2},$$

which is the moment characterization of sub-Weibull variable.

- **Extended envelope property** implies  $\|h^{(\ell)}\|_k \asymp \|g^{(\ell)}\|_k$
- **Base step:**  $g \sim N(0, \sigma^2)$ ,  $\|g\|_k \asymp \sqrt{k}$ . Thus,  $\|h\|_k = \|\phi(g)\|_k \asymp \|g\|_k \asymp \sqrt{k}$ .

## Penalized estimation

**Regularized problem:**

$$\min_{\mathbf{W}} \mathcal{R}(\mathbf{W}) + \lambda \mathcal{L}(\mathbf{W}), \quad (1)$$

where  $\mathcal{R}(\mathbf{W})$  is a **loss function**,  $\lambda \mathcal{L}(\mathbf{W})$  is a **penalty**,  $\lambda > 0$ .

For Bayesian models with prior distribution  $\pi(\mathbf{W})$ , the **maximum a posteriori** (MAP) solves (1) with:

$$\mathcal{L}(\mathbf{W}) \propto -\log \pi(\mathbf{W})$$

## Sparsity interpretation

**MAP on weights is L2-reg.**

Independent Gaussian prior

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

is equivalent to the weight decay penalty with *negative log-prior*:

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

**MAP on units induces sparsity**

The joint prior distribution for all the units can be expressed by Sklar's representation theorem as

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

where  $C$  is the copula of  $\mathbf{U}$  (characterizes all the dependence between the units),  $F$  is its cumulative distribution function. The penalty is the *negative log-prior*:

$$\mathcal{L}(\mathbf{U}) \approx \|\mathbf{U}^{(1)}\|_2^2 + \dots + \|\mathbf{U}^{(L)}\|_{2/L}^{2/L} - \log C(F(\mathbf{U})).$$

| Layer  | $\mathbf{W}$ -penalty                        | $\mathbf{U}$ -penalty   |
|--------|--|---|
| 1      | $\ \mathbf{W}^{(1)}\ _2^2, \mathcal{L}^2$    | $\ \mathbf{U}^{(1)}\ _2^2, \mathcal{L}^2$                         |
| 2      | $\ \mathbf{W}^{(2)}\ _2^2, \mathcal{L}^2$    | $\ \mathbf{U}^{(2)}\ , \mathcal{L}^1$                             |
| $\ell$ | $\ \mathbf{W}^{(\ell)}\ _2^2, \mathcal{L}^2$ | $\ \mathbf{U}^{(\ell)}\ _{2/\ell}^{2/\ell}, \mathcal{L}^{2/\ell}$ |

## Conclusion

We prove that the marginal prior unit distributions are heavier-tailed as depth increases. We further interpret this finding, showing that the units tend to be more sparsely represented as layers become deeper. This result provides new theoretical insight on deep Bayesian neural networks, underpinning their natural shrinkage properties and practical potential.

## References

Vladimirova, M., Arbel, J., and Mesejo, P. (2018). Bayesian neural networks increasingly sparsify their units with depth. *arXiv preprint arXiv:1810.05193*.