

---

# Supplementary Materials: Dropout as a Structured Shrinkage Prior

---

Eric Nalisnick<sup>1</sup> José Miguel Hernández-Lobato<sup>1,2,3</sup> Padhraic Smyth<sup>4</sup>

## 1. ADD as Multiplicative Noise and Equivalence to Stochastic Depth

Just as we did for ARD, we can re-formulate ADD in its expanded parametrization, revealing its equivalent MN form:

$$\mathbf{h}_{n,l} = f_l(\mathbf{h}_{n,l-1} \mathbf{W}_l) + \mathbf{h}_{n,l-1} \xrightarrow{\text{reparametrization}} f_l(\tau_l \mathbf{h}_{n,l-1} \mathbf{W}_l) + \mathbf{h}_{n,l-1} \quad (1)$$

where  $\tau_l$  is the ADD scale random variable. Comparing the equation above to the corresponding equation for ARD, we see that ADD is much less sampling intensive, requiring just one noise variable for each hidden layer.

Huang et al. (2016) proposed stochastic depth resnets by applying dropout to the resnet block, i.e.  $\mathbf{a}_l = \lambda_l F_l(\mathbf{h}_{l-1}) + \mathbf{h}_{l-1}$  where  $F(\cdot)$  denotes a whole resnet block and  $\lambda_l$  is a Bernoulli random variable. For simplicity, if we assume the resnet block consists of just one non-linear transformation, then applying Bernoulli ADD in MN form to the resnet architecture yields a similar expression, as seen in Equation 1. The only difference is that the Bernoulli variable is within  $f_l$  for ADD whereas it is outside the block for stochastic depth. However, if we assume  $f_l$  is the ReLU function, then the activation is scale equivariant. If the noise’s support is non-negative, then it is equivalent to move the noise variable outside the activation:  $\text{ReLU}(\tau_l \mathbf{h}_{l-1} \mathbf{W}_l) + \mathbf{h}_{l-1} = \tau_l \text{ReLU}(\mathbf{h}_{l-1} \mathbf{W}_l) + \mathbf{h}_{l-1}$  for  $\tau \geq 0$ . Thus, we can derive the previously proposed stochastic depth regularization as a special case of our ADD framework.

## 2. Parametrizations and their Effect on the Jensen’s Gap

Given the equivalence between training under multiplicative noise and MC integration of the scale prior, we wonder: are we working with the correct parametrization? A first glance suggests that sampling noise in the hierarchical parametrization, thereby injecting it into the prior, could provide a more stable MC estimator than one obtained by perturbing the NN. We investigate this hypothesis below through analysis of the Jensen’s gap of each estimator. It quickly becomes apparent that the hierarchical parametrization has undesirable properties when the hyper-prior is used as a regularizer. For all results, we assume  $\mathbb{E}[\xi^2] \leq 1$ , which is appropriate since we are interested in the settings in which the prior induces strong shrinkage, i.e.  $\xi^2 \approx 0$ . The precise quantity we wish to analyze is  $\mathcal{J}_{\text{GAP}} = \log p(\mathbf{y}, \{\mathbf{W}_l\}_{l=1}^{L+1} | \mathbf{X}) - \mathbb{E}_{p(\xi)} [\log p(\mathbf{y}, \{\mathbf{W}_l\}_{l=1}^{L+1} | \mathbf{X}, \{\Xi_l\}_{l=1}^L)]$  where the first term is the model joint with the scale integrated out and the second is the lower bound induced by invoking Jensen’s inequality to extract the expectation.

We first consider the model in its expanded parametrization. Denote the Jensen’s gap in this setting as  $\mathcal{J}_{\text{EP-GAP}}$ . The following proposition reveals the gap’s series representation.

**Proposition 2.1.** *Let the NN that parametrizes the likelihood be a 2-Lipschitz function, let  $k$  denote an integer, let  $\mu_\xi = \mathbb{E}[\xi]$ , and let  $\text{Var}[\xi]$  denote the variance of the scale (noise) distribution. The Jensen’s gap of the expanded parametrization is:*

$$\begin{aligned} \mathcal{J}_{\text{EP-GAP}} &= \frac{1}{2} \text{Var}[\xi] \left\| \nabla_{\mu_\xi} \log p(\mathbf{y} | \mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\Xi_l]\}_{l=1}^L) \right\|_2^2 - \sum_{k=2}^{\infty} \frac{(-1)^k}{k} \left[ \frac{1}{2} \text{Var}[\xi] \Omega \right]^k \\ &= \frac{1}{2} \text{Var}[\xi] \left\| \nabla_{\mu_\xi} \log p(\mathbf{y} | \mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\Xi_l]\}_{l=1}^L) \right\|_2^2 - \frac{1}{8} \text{Var}^2[\xi] \Omega^2 + \mathcal{O}(\text{Var}^3[\xi]) \end{aligned}$$

---

<sup>1</sup>Department of Engineering, University of Cambridge, Cambridge, United Kingdom <sup>2</sup>Microsoft Research, Cambridge, United Kingdom <sup>3</sup>Alan Turing Institute <sup>4</sup>Department of Computer Science, University of California, Irvine, United States of America. Correspondence to: Eric Nalisnick <e.nalisnick@eng.cam.ac.uk>.

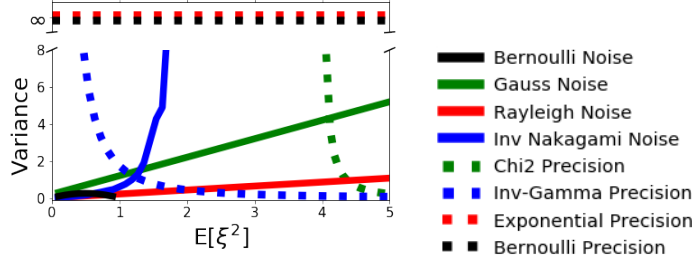


Figure 1. Noise Variances. The plot shows  $\text{Var}[\xi]$  in solid lines and  $\text{Var}[\xi^{-2}]$  in dashed lines for the noise distributions in Table 1 of the main text.

$$\text{where } \Omega = \left\| \nabla_{\mu_\xi} \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\Xi_l]\}_{l=1}^L) \right\|_2^2 + \text{Tr} \left\{ \nabla_{\mu_\xi}^2 \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\Xi_l]\}_{l=1}^L) \right\}$$

The proof is given in the supplementary materials. We see that the gap is a function of the variance of the scale and the norm of the log likelihood’s gradient w.r.t. the scales. This latter quantity could be large. However, as  $\mathbb{E}[\xi] \rightarrow 0$  so does  $\text{Var}[\xi] \rightarrow 0$ . We plot  $\mathbb{E}[\xi^2]$  vs  $\text{Var}[\xi]$  in Figure 1 (solid lines) for the noise distributions given in Table 1. Not only is  $\text{Var}[\xi]$  near zero when the prior is encouraging strong shrinkage but we also see favorable scaling for all but the inverse Nakagami as  $\mathbb{E}[\xi^2]$  grows.

We now move on to the hierarchical parametrization in which the random scale is left in the Gaussian prior. Denote the Jensen’s gap as  $\mathcal{J}_{\text{HP-GAP}}$ , and we give its series representation below.

**Proposition 2.2.** *Let  $p(\mathbf{W}|\Xi^{-2})$  be a factorized, zero-mean Gaussian prior with  $\xi^{-2}$  denoting its precision. The Jensen’s gap in the hierarchical parametrization is:*

$$\mathcal{J}_{\text{HP-GAP}} = \frac{L}{2} \text{Var}[\xi^{-2}] \left\| \nabla_{\mu_{\xi^{-2}}} \log p(\mathbf{W}|\mathbb{E}[\Xi^{-2}]) \right\|_2^2 - L \sum_{k=2}^{\infty} \frac{(-1)^k}{k} \left[ \frac{1}{2} \text{Var}[\xi^{-2}] \Psi \right]^k$$

$$\text{where } \Psi = \left\| \nabla_{\mu_{\xi^{-2}}} \log p(\mathbf{W}|\mathbb{E}[\Xi^{-2}]) \right\|_2^2 - \frac{1}{2\mathbb{E}[\xi^{-2}]}$$

The derivation is provided in the supplementary materials. We see that the gap depends on the variance of  $1/\xi^2$  and the norm of the prior’s gradient w.r.t. the precision. While the latter term is likely manageable, the former term is not. In Figure 1 we plot  $\mathbb{E}[\xi^2]$  vs  $\text{Var}[\xi^{-2}]$  (dashed lines) for the same scale distributions, and we see that either they asymptote as  $\mathbb{E}[\xi^2] \rightarrow 0$ —which occurs for Gaussian noise (green) at  $\sim 4$  and for inverse Nakagami noise (blue) at  $\sim 0.5$ —or are infinite for the whole range—Bernoulli (black) and Rayleigh (red). Thus, the series will not converge when the shrinkage is strong, only when  $\mathbb{E}[\xi^2]$  is well above one.

In summary, Proposition 2.1 shows that injecting noise into the likelihood relaxes the MAP estimate as a function of  $\text{Var}[\xi]$  whereas Proposition 2.2 reveals noise in the Gaussian prior has a  $\text{Var}[\xi^{-2}]$ -order gap. When  $\mathbb{E}[\xi^2]$  is near zero,  $\text{Var}[\xi^{-2}]$  explodes for all noise distributions considered (Figure 1), resulting in an impractical objective.

### Proposition 2.1: Expanded Parametrization Gap

$$\begin{aligned} \mathcal{J}_{\text{EP-GAP}} &= \log p(\mathbf{y}, \{\mathbf{W}_l\}_{l=1}^{L+1}|\mathbf{X}) - \mathbb{E}_{p(\xi)} \left[ \log p(\mathbf{y}, \{\mathbf{W}_l\}_{l=1}^{L+1}|\mathbf{X}, \{\Xi_l\}_{l=1}^L) \right] \\ &= \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}) + \log p(\{\mathbf{W}_l\}_{l=1}^{L+1}) - \mathbb{E}_{p(\xi)} \left[ \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\Xi_l\}_{l=1}^L) \right] - \log p(\{\mathbf{W}_l\}_{l=1}^{L+1}) \\ &= \log \mathbb{E}_{p(\xi)} \left[ p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\Xi_l\}_{l=1}^L) \right] - \mathbb{E}_{p(\xi)} \left[ \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\Xi_l\}_{l=1}^L) \right]. \end{aligned}$$

Expanding  $\mathbb{E}_{p(\xi)} \left[ p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\Xi_l\}_{l=1}^L) \right]$  around  $\mathbb{E}[\Xi]$ , we have

$$\mathbb{E}_{p(\xi)} \left[ p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\Xi_l\}_{l=1}^L) \right] \approx p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\Xi_l]\}_{l=1}^L) + \frac{1}{2} \text{Var}[\xi] \text{Tr} \left\{ \nabla_{\mu_{\xi}}^2 \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\Xi_l]\}_{l=1}^L) \right\},$$

and dividing through by  $p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\boldsymbol{\Xi}_l]\}_{l=1}^L)$  inside the log yields:

$$\begin{aligned} \log \mathbb{E}_{p(\xi)} [p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\boldsymbol{\Xi}_l\}_{l=1}^L)] &\approx \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\boldsymbol{\Xi}_l]\}_{l=1}^L) + \\ &\log \left\{ 1 + \frac{\frac{1}{2} \text{Var}[\xi] \text{Tr} \left\{ \nabla_{\mu_{xi}}^2 p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\boldsymbol{\Xi}_l]\}_{l=1}^L) \right\}}{p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\boldsymbol{\Xi}_l]\}_{l=1}^L)} \right\}. \end{aligned}$$

Noticing that

$$\begin{aligned} \nabla_{\mu_{xi}}^2 p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\boldsymbol{\Xi}_l]\}_{l=1}^L) = \\ p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\boldsymbol{\Xi}_l]\}_{l=1}^L) \left[ \left( \nabla_{\mu_{xi}} \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\boldsymbol{\Xi}_l]\}_{l=1}^L) \right)^2 + \nabla_{\mu_{xi}}^2 \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\boldsymbol{\Xi}_l]\}_{l=1}^L) \right] \end{aligned}$$

allows the previous expression to simplify to

$$\log \mathbb{E}_{p(\xi)} [p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\boldsymbol{\Xi}_l\}_{l=1}^L)] \approx \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\boldsymbol{\Xi}_l]\}_{l=1}^L) + \log \left\{ 1 + \frac{1}{2} \text{Var}[\xi] \text{Tr} \{ \boldsymbol{\Omega} \} \right\}$$

$$\text{where } \boldsymbol{\Omega} = \left\| \nabla_{\mu_{xi}} \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\boldsymbol{\Xi}_l]\}_{l=1}^L) \right\|_2^2 + \text{Tr} \left\{ \nabla_{\mu_{xi}}^2 \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\boldsymbol{\Xi}_l]\}_{l=1}^L) \right\}.$$

Next expanding the second term similarly around  $\mathbb{E}[\boldsymbol{\Xi}]$ , we have

$$\begin{aligned} \mathbb{E}_{p(\xi)} [\log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\boldsymbol{\Xi}_l\}_{l=1}^L)] &\approx \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\boldsymbol{\Xi}_l]\}_{l=1}^L) \\ &+ \frac{1}{2} \text{Var}[\xi] \text{Tr} \left\{ \nabla_{\mu_{xi}}^2 \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\boldsymbol{\Xi}_l]\}_{l=1}^L) \right\}. \end{aligned}$$

Now we put both terms together while expanding the  $\log\{\cdot\}$  around 1:

$$\begin{aligned} \mathcal{J}_{\text{EP-GAP}} &\approx \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\boldsymbol{\Xi}_l]\}_{l=1}^L) + \log 1 - \sum_{k=1}^{\infty} \frac{(-1)^k}{k} \left[ \frac{1}{2} \text{Var}[\xi] \boldsymbol{\Omega} \right]^k \\ &\quad - \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\boldsymbol{\Xi}_l]\}_{l=1}^L) - \frac{1}{2} \text{Var}[\xi] \text{Tr} \left\{ \nabla_{\mu_{xi}}^2 \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\boldsymbol{\Xi}_l]\}_{l=1}^L) \right\} \\ &= - \sum_{k=1}^{\infty} \frac{(-1)^k}{k} \left[ \frac{1}{2} \text{Var}[\xi] \boldsymbol{\Omega} \right]^k - \frac{1}{2} \text{Var}[\xi] \text{Tr} \left\{ \nabla_{\mu_{xi}}^2 \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\boldsymbol{\Xi}_l]\}_{l=1}^L) \right\} \\ &= \frac{1}{2} \text{Var}[\xi] \boldsymbol{\Omega} - \frac{1}{2} \text{Var}[\xi] \text{Tr} \left\{ \nabla_{\mu_{xi}}^2 \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\boldsymbol{\Xi}_l]\}_{l=1}^L) \right\} - \sum_{k=2}^{\infty} \frac{(-1)^k}{k} \left[ \frac{1}{2} \text{Var}[\xi] \boldsymbol{\Omega} \right]^k \\ &= \frac{1}{2} \text{Var}[\xi] \left\| \nabla_{\mu_{xi}} \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\boldsymbol{\Xi}_l]\}_{l=1}^L) \right\|_2^2 - \sum_{k=2}^{\infty} \frac{(-1)^k}{k} \left[ \frac{1}{2} \text{Var}[\xi] \boldsymbol{\Omega} \right]^k \\ &= \frac{1}{2} \text{Var}[\xi] \left\| \nabla_{\mu_{xi}} \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}, \{\mathbb{E}[\boldsymbol{\Xi}_l]\}_{l=1}^L) \right\|_2^2 - \frac{1}{8} \text{Var}^2[\xi] \boldsymbol{\Omega}^2 + \mathcal{O}(\text{Var}^3[\xi]). \end{aligned}$$

We assume that  $\text{Var}[\xi] \rightarrow 0$  (i.e. strong shrinkage) and this gives an  $\mathcal{O}(\text{Var}^3[\xi])$  convergence.

### Proposition 2.2: Hierarchical Parametrization Gap

$$\begin{aligned} \mathcal{J}_{\text{HP-GAP}} &= \log p(\mathbf{y}, \{\mathbf{W}_l\}_{l=1}^{L+1}|\mathbf{X}) - \mathbb{E}_{p(\xi)} [\log p(\mathbf{y}, \{\mathbf{W}_l\}_{l=1}^{L+1}|\mathbf{X}, \{\boldsymbol{\Xi}_l\}_{l=1}^L)] \\ &= \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}) + \log \mathbb{E}_{p(\xi)} [p(\{\mathbf{W}_l\}_{l=1}^{L+1}|\{\boldsymbol{\Xi}_l\}_{l=1}^L)] - \log p(\mathbf{y}|\mathbf{X}, \{\mathbf{W}_l\}_{l=1}^{L+1}) - \mathbb{E}_{p(\xi)} [\log p(\{\mathbf{W}_l\}_{l=1}^{L+1}|\{\boldsymbol{\Xi}_l\}_{l=1}^L)] \\ &= \log \mathbb{E}_{p(\xi)} [p(\{\mathbf{W}_l\}_{l=1}^{L+1}|\{\boldsymbol{\Xi}_l\}_{l=1}^L)] - \mathbb{E}_{p(\xi)} [\log p(\{\mathbf{W}_l\}_{l=1}^{L+1}|\{\boldsymbol{\Xi}_l\}_{l=1}^L)] \\ &= L \log \mathbb{E}_{p(\xi)} [p(\mathbf{W}|\boldsymbol{\Xi})] - L \mathbb{E}_{p(\xi)} [\log p(\mathbf{W}|\boldsymbol{\Xi})]. \end{aligned}$$

Next we reparametrize so that  $\boldsymbol{\Xi}^{-2}$  denotes the Gaussian's precision. We then expand the first term around  $\mathbb{E}[\boldsymbol{\Xi}^{-2}]$ :

$$\mathbb{E}_{p(\xi)} [p(\mathbf{W}|\boldsymbol{\Xi}^{-2})] \approx p(\mathbf{W}|\mathbb{E}[\boldsymbol{\Xi}^{-2}]) + \frac{1}{2} \text{Var}[\xi^{-2}] \text{Tr} \left\{ \nabla_{\mu_{\xi^{-2}}}^2 p(\mathbf{W}|\mathbb{E}[\boldsymbol{\Xi}^{-2}]) \right\}.$$

Dividing inside the log then allows us to write:

$$L \log \mathbb{E}_{p(\xi)} [p(\mathbf{W}|\Xi^{-2})] \approx L \log p(\mathbf{W}|\mathbb{E}[\Xi^{-2}]) + L \log \left\{ 1 + \frac{\frac{1}{2} \text{Var}[\xi^{-2}] \text{Tr} \left\{ \nabla_{\mu_{\xi^{-2}}}^2 p(\mathbf{W}|\mathbb{E}[\Xi^{-2}]) \right\}}{p(\mathbf{W}|\mathbb{E}[\Xi^{-2}])} \right\}.$$

Noticing that

$$\nabla_{\mu_{\xi^{-2}}}^2 p(\mathbf{W}|\mathbb{E}[\Xi^{-2}]) = p(\mathbf{W}|\mathbb{E}[\Xi^{-2}]) \left[ \left( \nabla_{\mu_{\xi^{-2}}} \log p(\mathbf{W}|\mathbb{E}[\Xi^{-2}]) \right)^2 + \nabla_{\mu_{\xi^{-2}}}^2 \log p(\mathbf{W}|\mathbb{E}[\Xi^{-2}]) \right]$$

allows the previous expression to simplify to

$$L \log \mathbb{E}_{p(\xi)} [p(\mathbf{W}|\Xi^{-2})] \approx L \log p(\mathbf{W}|\mathbb{E}[\Xi^{-2}]) + L \log \left\{ 1 + \frac{1}{2} \text{Var}[\xi^{-2}] \Psi \right\}$$

where

$$\Psi = \|\nabla_{\mu_{\xi^{-2}}} \log p(\mathbf{W}|\mathbb{E}[\Xi^{-2}])\|_2^2 + \text{Tr} \left\{ \nabla_{\mu_{\xi^{-2}}}^2 \log p(\mathbf{W}|\mathbb{E}[\Xi^{-2}]) \right\}.$$

Next we can write out the second term as

$$\begin{aligned} L \mathbb{E}_{p(\xi)} [\log p(\mathbf{W}|\Xi^{-2})] &= \frac{-L}{2\sigma_0^2} \mathbb{E}[\xi^{-2}] \|\mathbf{W}\|_2^2 + \frac{L}{2} \mathbb{E}[\log \xi^{-2}] - \frac{L}{2} \log 2\pi \\ &\approx \frac{-L}{2\sigma_0^2} \mathbb{E}[\xi^{-2}] \|\mathbf{W}\|_2^2 + \frac{L}{2} \log \mathbb{E}[\xi^{-2}] - \frac{L}{2} \log 2\pi + \frac{L}{4} \text{Var}[\xi^{-2}] \frac{-1}{\mathbb{E}[\xi^{-2}]^2} \\ &= L \log p(\mathbf{W}|\mathbb{E}[\Xi^{-2}]) - \text{Var}[\xi^{-2}] \frac{L}{4\mathbb{E}[\xi^{-2}]^2}. \end{aligned}$$

Rejoining the expressions, again expanding the log around 1, we have:

$$\begin{aligned} \mathcal{J}_{\text{HP-GAP}} &\approx L \log p(\mathbf{W}|\mathbb{E}[\Xi^{-2}]) + L \log 1 - L \sum_{k=1}^{\infty} \frac{(-1)^k}{k} \left[ \frac{1}{2} \text{Var}[\xi^{-2}] \Psi \right]^k - L \log p(\mathbf{W}|\mathbb{E}[\Xi^{-2}]) + \text{Var}[\xi^{-2}] \frac{L}{4\mathbb{E}[\xi^{-2}]^2} \\ &= -L \sum_{k=1}^{\infty} \frac{(-1)^k}{k} \left[ \frac{1}{2} \text{Var}[\xi^{-2}] \Psi \right]^k + \text{Var}[\xi^{-2}] \frac{L}{4\mathbb{E}[\xi^{-2}]^2} \\ &= \frac{L}{2} \text{Var}[\xi^{-2}] \Psi + \frac{L}{2} \text{Var}[\xi^{-2}] \frac{1}{2\mathbb{E}[\xi^{-2}]^2} - L \sum_{k=2}^{\infty} \frac{(-1)^k}{k} \left[ \frac{1}{2} \text{Var}[\xi^{-2}] \Psi \right]^k \\ &= \frac{L}{2} \text{Var}[\xi^{-2}] \|\nabla_{\mu_{\xi^{-2}}} \log p(\mathbf{W}|\mathbb{E}[\Xi^{-2}])\|_2^2 - L \sum_{k=2}^{\infty} \frac{(-1)^k}{k} \left[ \frac{1}{2} \text{Var}[\xi^{-2}] \Psi \right]^k. \end{aligned}$$

### 3. Variational EM Updates

#### 3.1. ADD EM Updates

Below we derive the EM updates for the various hyper-priors considered. For all expressions, we assume the variational distribution on the weights factorizes, i.e.  $\mathbf{N}(\mathbf{W}; \phi) = \prod_i \prod_j \mathbf{N}(w_{i,j}; \phi_{i,j})$  where  $i$  denotes row and  $j$  column indices.

##### 3.1.1. INVERSE GAMMA

We begin with the ELBO terms that depend on the variational variance  $\tau$ :

$$\begin{aligned} \mathcal{J}_{\text{ELBO}}(\phi, \bar{\tau}_l) &= -\text{KL} [\mathbf{N}(\mathbf{W}; \phi) | | \mathbf{N}(\mathbf{W} | \bar{\tau}_l)] + \log p(\bar{\tau}_l) \\ &= \frac{-1}{2} \sum_i \sum_j \left[ \log \frac{\bar{\tau}_l}{\sigma_{i,j}^2} + \frac{\sigma_{i,j}^2}{\bar{\tau}_l} + \frac{\mu_{i,j}^2}{\bar{\tau}_l} - 1 \right] + \log \Gamma^{-1}(\bar{\tau}_l; \alpha, \beta) \\ &= \frac{-1}{2} \left[ \sum_i \sum_j \log \frac{\bar{\tau}_l}{\sigma_{i,j}^2} + \sum_i \sum_j \frac{\sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\tau}_l} - D \right] - \frac{\beta}{\bar{\tau}_l} - (\alpha + 1) \log \bar{\tau}_l + \log \frac{\beta^\alpha}{\Gamma(\alpha)}. \end{aligned}$$

Differentiating w.r.t.  $\bar{\tau}_l$  and setting to zero, we have

$$\begin{aligned}
 0 &= \frac{\partial}{\partial \bar{\tau}_l} \mathcal{J}_{\text{ELBO}}(\phi, \bar{\tau}_l) \\
 &= \frac{-1}{2} \left[ \sum_i \sum_j \frac{1}{\bar{\tau}_l} - \frac{\sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\tau}_l^2} \right] + \frac{\beta}{\bar{\tau}_l^2} - \frac{\alpha + 1}{\bar{\tau}_l} \\
 &= \frac{D}{\bar{\tau}_l} - \frac{\sum_i \sum_j \sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\tau}_l^2} - \frac{2\beta}{\bar{\tau}_l^2} + \frac{2\alpha + 2}{\bar{\tau}_l} \\
 \frac{2\beta + \sum_i \sum_j \sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\tau}_l^2} &= \frac{D + 2\alpha + 2}{\bar{\tau}_l} \\
 \bar{\tau}_l^* &= \frac{2\beta + \sum_i \sum_j \sigma_{i,j}^2 + \mu_{i,j}^2}{D + 2\alpha + 2}
 \end{aligned}$$

where  $D$  are the number of dimensions (i.e. parameters) in  $W$ .

### 3.1.2. HALF-CAUCHY

Again starting with the ELBO, we have

$$\begin{aligned}
 0 &= \frac{\partial}{\partial \bar{\tau}_l} \mathcal{J}_{\text{ELBO}}(\phi, \bar{\tau}_l) \\
 &= \frac{-1}{2} \frac{\partial}{\partial \bar{\tau}_l} \sum_i \sum_j \left[ \log \frac{\bar{\tau}_l}{\sigma_{i,j}^2} + \frac{\sigma_{i,j}^2}{\bar{\tau}_l} + \frac{\mu_{i,j}^2}{\bar{\tau}_l} - 1 \right] + \frac{\partial}{\partial \bar{\tau}_l} \log C^+(\sqrt{\bar{\tau}_l}; 0, \eta) \\
 &= \frac{-1}{2} \left[ \sum_i \sum_j \frac{1}{\bar{\tau}_l} - \frac{\sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\tau}_l^2} \right] + \frac{\partial}{\partial \bar{\tau}_l} \log \frac{2}{\pi \eta (1 + \bar{\tau}_l / \eta^2)} \\
 &= \frac{D}{\bar{\tau}_l} - \frac{\sum_i \sum_j \sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\tau}_l^2} + \frac{2}{\eta^2 + \bar{\tau}_l} \\
 &= \frac{2\bar{\tau}_l}{\eta^2 + \bar{\tau}_l} - \frac{\sum_i \sum_j \sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\tau}_l} + D.
 \end{aligned}$$

Solving<sup>1</sup> the above equation for  $\bar{\tau}_l$  gives for a positive solution:

$$\bar{\tau}_l^* = \frac{M - \eta^2 D + \sqrt{M^2 + (2D + 8)\eta^2 M + \eta^4 D^2}}{2D + 4} \quad \text{where } M = \sum_i \sum_j \sigma_{i,j}^2 + \mu_{i,j}^2,$$

$D$  is again the dimensionality, and  $\eta$  is the half-Cauchy's scale.

<sup>1</sup>We plugged the equation into Wolfram Alpha.

## 3.1.3. LOG-UNIFORM

Again beginning with ELBO, we have

$$\begin{aligned}
 0 &= \frac{\partial}{\partial \bar{\tau}_l} \mathcal{J}_{\text{ELBO}}(\phi, \bar{\tau}_l) \\
 &= \frac{-1}{2} \left[ \sum_i \sum_j \frac{1}{\bar{\tau}_l} - \frac{\sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\tau}_l^2} \right] + \frac{\partial}{\partial \bar{\tau}_l} \log \frac{c}{\bar{\tau}_l} \\
 &= \frac{D}{\bar{\tau}_l} - \frac{\sum_i \sum_j \sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\tau}_l^2} + \frac{2}{\bar{\tau}_l} \\
 \frac{\sum_i \sum_j \sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\tau}_l^2} &= \frac{D+2}{\bar{\tau}_l} \\
 \bar{\tau}_l^* &= \frac{\sum_i \sum_j \sigma_{i,j}^2 + \mu_{i,j}^2}{D+2}.
 \end{aligned}$$

Notice that this update is the same as the inverse Gamma's as  $\alpha \rightarrow 0$  and  $\beta \rightarrow 0$ .

## 3.2. ARD-ADD EM Updates

For approximate posterior inference for the joint ARD-ADD prior, we assume the posterior approximation

$$p(\mathbf{W}, \lambda, \tau | \mathbf{y}, \mathbf{X}) \approx q(\mathbf{W}; \phi) q(\lambda) q(\tau) = \text{N}(\mathbf{W}; \boldsymbol{\mu}_\phi, \text{diag}\{\boldsymbol{\Sigma}_\phi\}) \delta[\bar{\lambda}_i] \delta[\bar{\tau}_l]$$

where  $\phi = \{\boldsymbol{\mu}_\phi, \text{diag}\{\boldsymbol{\Sigma}_\phi\}\}$ ,  $\bar{\lambda}_i$ , and  $\bar{\tau}_l$  are the variational parameters. The ELBO for this approximation is

$$\begin{aligned}
 \log p(\mathbf{y} | \mathbf{X}) &\geq \\
 \mathbb{E}_{q(\mathbf{W})} [\log p(\mathbf{y} | \mathbf{X}, \mathbf{W})] &- \mathbb{E}_{q(\lambda)} \mathbb{E}_{q(\tau)} \text{KL} [q(\mathbf{W}; \phi) || p(\mathbf{W} | \lambda, \tau)] - \text{KL} [q(\lambda) || p(\lambda)] - \text{KL} [q(\tau) || p(\tau)] \\
 &= \mathbb{E}_{\text{N}(\mathbf{W})} [\log p(\mathbf{y} | \mathbf{X}, \mathbf{W})] - \text{KL} [\text{N}(\mathbf{W}; \phi) || \text{N}(\mathbf{W} | \bar{\lambda}_i, \bar{\tau}_l)] - \text{KL} [\delta[\bar{\lambda}_i] || p(\lambda)] - \text{KL} [\delta[\bar{\tau}_l] || p(\tau)] \\
 &= \mathbb{E}_{\text{N}(\mathbf{W})} [\log p(\mathbf{y} | \mathbf{X}, \mathbf{W})] - \text{KL} [\text{N}(\mathbf{W}; \phi) || \text{N}(\mathbf{W} | \bar{\lambda}_i, \bar{\tau}_l)] + \log p(\bar{\lambda}_i) + \log p(\bar{\tau}_l) + \mathcal{C}
 \end{aligned} \tag{2}$$

where  $\mathcal{C} = \mathbb{H}[\delta[\bar{\lambda}_i]] + \mathbb{H}[\delta[\bar{\tau}_l]]$  is again a constant. Next we must find solutions for both  $\bar{\lambda}_i$  and  $\bar{\tau}_l$ . Again we can differentiate the ELBO and set to zero:

$$\begin{aligned}
 \frac{\partial}{\partial \bar{\lambda}_i} \mathcal{J}_{\text{ELBO}}(\phi, \bar{\lambda}_i, \bar{\tau}_l) &= -\frac{\partial}{\partial \bar{\lambda}_i} \text{KL} [\text{N}(\mathbf{W}; \phi) || \text{N}(\mathbf{W} | \bar{\lambda}_i, \bar{\tau}_l)] + \frac{\partial}{\partial \bar{\lambda}_i} \log p(\bar{\lambda}_i) = 0, \\
 \frac{\partial}{\partial \bar{\tau}_l} \mathcal{J}_{\text{ELBO}}(\phi, \bar{\lambda}_i, \bar{\tau}_l) &= -\frac{\partial}{\partial \bar{\tau}_l} \text{KL} [\text{N}(\mathbf{W}; \phi) || \text{N}(\mathbf{W} | \bar{\lambda}_i, \bar{\tau}_l)] + \frac{\partial}{\partial \bar{\tau}_l} \log p(\bar{\tau}_l) = 0.
 \end{aligned} \tag{3}$$

We denote the solutions as  $\bar{\lambda}_i^*$  and  $\bar{\tau}_l^*$ , deriving them for the three hyper-priors below. Again, we update  $q(\mathbf{W}; \phi)$  via gradient ascent:

$$\frac{\partial}{\partial \phi} \mathcal{J}_{\text{ELBO}}(\phi, \bar{\lambda}_i^*, \bar{\tau}_l^*) = \frac{\partial}{\partial \phi} \mathbb{E}_{\text{N}(\mathbf{W}; \phi)} [\log p(\mathbf{y} | \mathbf{X}, \mathbf{W})] - \frac{\partial}{\partial \phi} \text{KL} [\text{N}(\mathbf{W}; \phi) || \text{N}(\mathbf{W} | \bar{\lambda}_i^*, \bar{\tau}_l^*)]. \tag{4}$$

## 3.2.1. INVERSE GAMMA

Starting with the first line of Equation 3, we have

$$\begin{aligned}
 0 &= \frac{\partial}{\partial \bar{\lambda}_i} \mathcal{J}_{\text{ELBO}}(\phi, \bar{\lambda}_i, \bar{\tau}_l) \\
 &= \frac{-1}{2} \frac{\partial}{\partial \bar{\lambda}_i} \sum_i \sum_j \left[ \log \frac{\bar{\lambda}_i \bar{\tau}_l}{\sigma_{i,j}^2} + \frac{\sigma_{i,j}^2}{\bar{\lambda}_i \bar{\tau}_l} + \frac{\mu_{i,j}^2}{\bar{\lambda}_i \bar{\tau}_l} - 1 \right] + \frac{\partial}{\partial \bar{\lambda}_i} \log \Gamma^{-1}(\bar{\lambda}_i; \alpha, \beta) \\
 &= \frac{-1}{2} \left[ \sum_j \frac{1}{\bar{\lambda}_i} - \frac{\sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\lambda}_i^2 \bar{\tau}_l} \right] + \frac{\beta}{\bar{\lambda}_i^2} - \frac{\alpha + 1}{\bar{\lambda}_i} \\
 &= \frac{D_i}{\bar{\lambda}_i} - \frac{\sum_j \sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\lambda}_i^2 \bar{\tau}_l} - \frac{2\beta}{\bar{\lambda}_i^2} + \frac{2\alpha + 2}{\bar{\lambda}_i} \\
 \frac{2\beta \bar{\tau}_l + \sum_j \sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\lambda}_i^2 \bar{\tau}_l} &= \frac{D_i + 2\alpha + 2}{\bar{\lambda}_i} \\
 \bar{\lambda}_i^* &= \frac{2\beta \bar{\tau}_l + \sum_j \sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\tau}_l (D_i + 2\alpha + 2)}
 \end{aligned}$$

where  $D_i$  is the number of parameters in the  $i$ th row (i.e.  $\lambda$ 's corresponding row) of the weight matrix  $\mathbf{W}$ . Furthermore, notice that the sum over posterior parameters is across only columns ( $j$  index).

Moving on the ADD parameter, we begin with the second line of Equation 3:

$$\begin{aligned}
 0 &= \frac{\partial}{\partial \bar{\tau}_l} \mathcal{J}_{\text{ELBO}}(\phi, \bar{\lambda}_i, \bar{\tau}_l) \\
 &= \frac{-1}{2} \frac{\partial}{\partial \bar{\tau}_l} \sum_i \sum_j \left[ \log \frac{\bar{\lambda}_i \bar{\tau}_l}{\sigma_{i,j}^2} + \frac{\sigma_{i,j}^2}{\bar{\lambda}_i \bar{\tau}_l} + \frac{\mu_{i,j}^2}{\bar{\lambda}_i \bar{\tau}_l} - 1 \right] + \frac{\partial}{\partial \bar{\tau}_l} \log \Gamma^{-1}(\bar{\tau}_l; \alpha, \beta) \\
 &= \frac{-1}{2} \left[ \sum_i \sum_j \frac{1}{\bar{\tau}_l} - \frac{\sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\lambda}_i \bar{\tau}_l^2} \right] + \frac{\beta}{\bar{\tau}_l^2} - \frac{\alpha + 1}{\bar{\tau}_l} \\
 &= \frac{D}{\bar{\tau}_l} - \frac{\sum_i \bar{\lambda}_i^{-1} \sum_j \sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\tau}_l^2} - \frac{2\beta}{\bar{\tau}_l^2} + \frac{2\alpha + 2}{\bar{\tau}_l} \\
 \frac{2\beta + \sum_i \bar{\lambda}_i^{-1} \sum_j \sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\tau}_l^2} &= \frac{D + 2\alpha + 2}{\bar{\tau}_l} \\
 \bar{\tau}_l^* &= \frac{2\beta + \sum_i \bar{\lambda}_i^{-1} \sum_j \sigma_{i,j}^2 + \mu_{i,j}^2}{D + 2\alpha + 2}.
 \end{aligned}$$

## 3.2.2. HALF-CAUCHY

Starting with the first line of Equation 3, we have

$$\begin{aligned}
 0 &= \frac{\partial}{\partial \bar{\lambda}_i} \mathcal{J}_{\text{ELBO}}(\phi, \bar{\lambda}_i, \bar{\tau}_l) \\
 &= \frac{-1}{2} \frac{\partial}{\partial \bar{\lambda}_i} \sum_i \sum_j \left[ \log \frac{\bar{\lambda}_i \bar{\tau}_l}{\sigma_{i,j}^2} + \frac{\sigma_{i,j}^2}{\bar{\lambda}_i \bar{\tau}_l} + \frac{\mu_{i,j}^2}{\bar{\lambda}_i \bar{\tau}_l} - 1 \right] + \frac{\partial}{\partial \bar{\lambda}_i} \log C^+(\sqrt{\bar{\lambda}_i}; 0, \eta) \\
 &= \sum_j \frac{1}{\bar{\lambda}_i} - \frac{\sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\lambda}_i^2 \bar{\tau}_l} + \frac{2}{\eta^2 + \bar{\lambda}_i} \\
 &= \frac{2\bar{\lambda}_i}{\eta^2 + \bar{\lambda}_i} - \frac{\bar{\tau}_l^{-1} \sum_j \sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\lambda}_i} + D_i
 \end{aligned}$$

where  $D_i$  is the number of parameters in the  $i$ th row. This final expression is of the same form we solved for ADD and therefore has the same solution with appropriately adjusted constants:

$$\bar{\lambda}_i^* = \frac{M_i - \eta^2 D_i + \sqrt{M_i^2 + (2D_i + 8)\eta^2 M_i + \eta^4 D_i^2}}{2D_i + 4} \quad \text{where } M_i = \bar{\tau}_l^{-1} \sum_j \sigma_{i,j}^2 + \mu_{i,j}^2.$$

Moving on the ADD parameter, we consider the second line of Equation 3:

$$\begin{aligned} 0 &= \frac{\partial}{\partial \bar{\tau}_l} \mathcal{J}_{\text{ELBO}}(\phi, \bar{\lambda}_i, \bar{\tau}_l) \\ &= \frac{-1}{2} \frac{\partial}{\partial \bar{\tau}_l} \sum_i \sum_j \left[ \log \frac{\bar{\lambda}_i \bar{\tau}_l}{\sigma_{i,j}^2} + \frac{\sigma_{i,j}^2}{\bar{\lambda}_i \bar{\tau}_l} + \frac{\mu_{i,j}^2}{\bar{\lambda}_i \bar{\tau}_l} - 1 \right] + \frac{\partial}{\partial \bar{\tau}_l} \log C^+(\sqrt{\bar{\tau}_l}; 0, \eta) \\ &= \sum_i \sum_j \frac{1}{\bar{\tau}_l} - \frac{\sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\lambda}_i \bar{\tau}_l^2} + \frac{2}{\eta^2 + \bar{\tau}_l} \\ &= \frac{D}{\bar{\tau}_l} - \frac{\sum_i \bar{\lambda}_i^{-1} \sum_j \sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\tau}_l^2} + \frac{2}{\eta^2 + \bar{\tau}_l} \\ &= \frac{2\bar{\tau}_l}{\eta^2 + \bar{\tau}_l} - \frac{\sum_i \bar{\lambda}_i^{-1} \sum_j \sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\tau}_l} + D. \end{aligned}$$

Again solving the same equation yields the solution:

$$\bar{\tau}_l^* = \frac{M_\lambda - \eta^2 D + \sqrt{M_\lambda^2 + (2D + 8)\eta^2 M_\lambda + \eta^4 D^2}}{2D + 4} \quad \text{where } M_\lambda = \sum_i \bar{\lambda}_i^{-1} \sum_j \sigma_{i,j}^2 + \mu_{i,j}^2.$$

### 3.2.3. LOG-UNIFORM

As mentioned earlier, the solution for the log-uniform distribution can be attained in the limit of the inverse Gamma. Plugging  $\alpha = 0$  and  $\beta = 0$  into the inverse gamma's solutions we obtain:

$$\bar{\lambda}_i^* = \frac{\sum_j \sigma_{i,j}^2 + \mu_{i,j}^2}{\bar{\tau}_l (D_i + 2)} \quad \text{and} \quad \bar{\tau}_l^* = \frac{\sum_i \bar{\lambda}_i^{-1} \sum_j \sigma_{i,j}^2 + \mu_{i,j}^2}{D + 2}.$$

## 4. Experimental Details

Below we provide more details on our experimental set-up. For all experiments, we used following settings, replicating the set-up of Gal & Ghahramani (2016). Optimization was done using Adam (Kingma & Ba, 2014) with a learning rate of .001 and other parameters kept at Tensorflow defaults. Training was done for 4500 epochs using batch sizes of 32 and 10 samples when calculating Monte Carlo expectations. The UCI regression data sets were divided into 20 randomized 90% – 10% train-test splits. Results are reported by training and evaluating on each split and reporting the average test metric. Test set evaluations were done by drawing 500 MC samples. The training split was standardized to zero mean and unit variance, and predictions were made by normalizing the test input based on these train set statistics and then applying the inverse transform to the prediction. RMSE and likelihood were calculated after applying the inverse transform.

**Importance Sampling Experiments** All results in Table 2 of the main text used a one-hidden-layer neural network with 50 hidden units and ReLU activations. We used the dropout rates and likelihood precisions reported in Table 1. These parameters were obtained from the validation set results reported by Gal & Ghahramani (2016) at: [https://github.com/yaringal/DropoutUncertaintyExps/tree/master/UCI\\_Datasets](https://github.com/yaringal/DropoutUncertaintyExps/tree/master/UCI_Datasets).

**Variational EM Experiments** The ARD, ADD, and ARD-ADD results in Table 3 of the main text used neural networks with two hidden layers with one skip connection between them. ReLU activations were again used exclusively. ARD was applied to the first and last weight matrices—since ADD can't be applied without a skip connection—and thus the



	Dropout Parameter ( $p$ )	Likelihood Precision
Boston	.050	0.20
Concrete	.005	0.075
Energy	.005	0.75
Kin8nm	.010	250.00
Power	.005	0.10
Wine	.050	3.50
Yacht	.005	0.75

Table 1. The table above reports the dropout rate and likelihood precision used for each experiment. We use the same parameters as Gal & Ghahramani (2016).

differences in performance are due to the prior’s influence on the hidden-to-hidden weights. We used *flipout* (Wen et al., 2018) to reduce the variance of the Monte Carlo expectations (10 samples). All scale priors were set to the same distribution, and we chose this distribution based on preliminary experiments on a subset of the Protein UCI data set (Dheeru & Karra Taniskidou, 2017). We tried the following priors:  $\Gamma^{-1}(0, 0)$  / log-uniform,  $\Gamma^{-1}(1, 1)$ ,  $\Gamma^{-1}(3, 3)$ ,  $\Gamma^{-1}(5, 5)$ ,  $\Gamma^{-1}(1, 3)$ ,  $\Gamma^{-1}(1, 5)$ ,  $\Gamma^{-1}(1, 10)$ ,  $C^+(0, .5)$ ,  $C^+(0, 1.)$ ,  $C^+(0, 2.5)$ ,  $C^+(0, 5)$ ,  $C^+(0, 10.)$ .  $\Gamma^{-1}(3, 3)$  performed best and we used it for all experiments. The fact that this one prior that was chosen on a different data set generalized so well speaks to the robustness of the EM algorithm.

## 5. Tensorflow Implementations

We have provided Python code to replicate the results presented in Tables 2 and 3 of the main text. We include the Yacht data set for the user’s convenience; the other data sets can be downloaded from Gal & Ghahramani (2016)’s GitHub repository<sup>2</sup> or from the UCI repository (Dheeru & Karra Taniskidou, 2017). Below we give implementation snippets for the paper’s core algorithmic contributions.

### 5.1. Implementation of Tail-Adaptive Importance Weights

```

1 log_likelihoods = # matrix of size Batch_Size x Number_of_MC_Samples
2
3 # Sort the samples
4 log_likelihoods_sorted, _ = tf.nn.top_k(log_likelihoods, k=n_samples, sorted=True)
5
6 # Use tf.range to provide ranks
7 gamma_vals = n_samples / (tf.expand_dims(tf.range(n_samples, dtype=tf.float32), 0) + 1.)
8
9 # Normalize and treat weight as a constant
10 normalized_weights = tf.stop_gradient(gamma_vals / tf.reduce_sum(gamma_vals, axis=1,
11     keep_dims=True))
12
13 # Apply weights
14 final_ell_term = tf.reduce_sum(normalized_weights * log_likelihoods_sorted,
15     reduction_indices=1, keep_dims=True)
16
17 # Return final expected log-likelihood term s.t. gradient will be w * d_ell
18 return tf.reduce_mean(final_ell_term)

```

### 5.2. EM Scale Updates

```

1 def ard_scale_update(prior_id, prior_params, mu, std):
2     param_count = mu.get_shape().as_list()[1]
3     moments = tf.reduce_sum(mu*mu + std*std, reduction_indices=1, keep_dims=True)
4
5     if prior_id == "inv_gamma":

```

<sup>2</sup>[https://github.com/yaringal/DropoutUncertaintyExps/tree/master/UCI\\_Datasets](https://github.com/yaringal/DropoutUncertaintyExps/tree/master/UCI_Datasets)

```

6     prior_opt_scale = tf.sqrt((2. * prior_params['param2'] + moments) / (param_count
7     + 2. * prior_params['param1'] + 2.))
8
9     elif prior_id == "half_cauchy":
10    prior_opt_scale = tf.sqrt((moments - prior_params['param1']**2 * param_count + tf.
11    sqrt(moments**2 + (2*param_count+8) * prior_params['param1']**2 * moments +
12    prior_params['param1']**4 * param_count**2)) / (2*param_count + 4))
13
14    # Log-uniform implemented by just passing alpha=beta=0 to inv Gamma
15
16    return prior_opt_scale * tf.ones_like(std)
17
18 def add_scale_update(prior_id, prior_params, mu, std):
19     param_count = mu.get_shape().as_list()[0] * mu.get_shape().as_list()[1]
20     moments = tf.reduce_sum(mu**2 + std**2)
21
22     if prior_id == "inv_gamma":
23         prior_opt_scale = tf.sqrt((2. * prior_params['param2'] + moments) / (param_count
24         + 2. * prior_params['param1'] + 2.))
25
26         elif prior_id == "half_cauchy":
27             prior_opt_scale = tf.sqrt((moments - prior_params['param1']**2 * param_count + tf.
28             sqrt(moments**2 + (2*param_count+8.) * prior_params['param1']**2 * moments +
29             prior_params['param1']**4 * param_count**2)) / (2*param_count + 4))
30
31         return prior_opt_scale * tf.ones_like(std)
32
33 def ard_add_scale_update(prior_id, prior_params, old_tau, mu, std):
34     param_count_cols = mu.get_shape().as_list()[1]
35     param_count_rows = mu.get_shape().as_list()[0]
36     moments = tf.reduce_sum(mu**2 + std**2, reduction_indices=1, keep_dims=True)
37
38     if prior_id == "inv_gamma":
39         opt_lambda = tf.sqrt((2. * prior_params['param2'] * old_tau + moments) / (old_tau
40         * 2. * prior_params['param1'] + old_tau * param_count_cols + old_tau * 2))
41         opt_tau = tf.sqrt((2. * prior_params['param2'] + tf.reduce_sum((1./opt_lambda) *
42         moments)) / (param_count_cols * param_count_rows + 2 + 2 * prior_params['param1']))
43
44     elif prior_id == "half_cauchy":
45         moments_tau = (1./old_tau) * moments
46         opt_lambda = tf.sqrt((moments_tau - prior_params['param1']**2 * param_count_cols +
47         tf.sqrt(moments_tau**2 + (2*param_count_cols+8) * prior_params['param1']**2 *
48         moments_tau + prior_params['param1']**4 * param_count_cols**2)) / (2*param_count_cols
49         + 4))
50
51         moments_lamb = tf.reduce_sum((1./opt_lambda) * moments)
52         param_count = param_count_rows * param_count_cols
53         opt_tau = tf.sqrt((moments_lamb - prior_params['param1']**2 * param_count + tf.
54         sqrt(moments_lamb**2 + (2*param_count+8) * prior_params['param1']**2 * moments_lamb +
55         prior_params['param1']**4 * param_count**2)) / (2*param_count + 4))
56
57         # Log-uniform implemented by just passing alpha=beta=0 to inv Gamma
58
59         return opt_lambda * opt_tau * tf.ones_like(std)

```

## References

- Dheeru, D. and Karra Taniskidou, E. UCI Machine Learning Repository, 2017. URL <http://archive.ics.uci.edu/ml>.
- Gal, Y. and Ghahramani, Z. Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning. In *Proceedings of the 33rd International Conference on Machine Learning (ICML)*, pp. 1050–1059, 2016.

## Supplementary Materials

---

Huang, G., Sun, Y., Liu, Z., Sedra, D., and Weinberger, K. Q. Deep Networks with Stochastic Depth. In *European Conference on Computer Vision (ECCV)*, pp. 646–661. Springer, 2016.

Kingma, D. and Ba, J. Adam: A Method for Stochastic Optimization. *International Conference on Learning Representations (ICLR)*, 2014.

Wen, Y., Vicol, P., Ba, J., Tran, D., and Grosse, R. Flipout: Efficient Pseudo-Independent Weight Perturbations on Mini-Batches. In *International Conference on Learning Representations*, 2018.