

Boosting Performance of Latent Binary Neural Networks With a Local Reparametrization Trick and Normalizing Flows

Lars Skaaret-Lund² **Aliaksandr Hubin**^{1,2,4} Geir Storvik^{1,3}

¹ Department of Mathematics, UiO

² BIAS, NMBU

³ SAMBA, NR

⁴ OUC

25.11.2022

Introduction. Issues with neural networks

- Neural networks (NNs) are flexible parametric models;
- But NNs overfit and do not provide uncertainty measures;
- Bayesian neural networks (BNNs) resolve that;
- BNNs are still heavily over-parameterized and uninterpretable;
- A solution is model uncertainties in BNNs;
- We develop latent binary BNN (LBBNN) for that.

The likelihood model for LBBNN

$$\mathbf{y}_i \sim f(\boldsymbol{\mu}_i, \phi), \quad i \in \{1, \dots, n\} \quad (1)$$

$$\boldsymbol{\mu}_i = \{x_{i1}^{(L)}, \dots, x_{ir}^{(L)}\}, \quad (2)$$

$$x_{ij}^{(l+1)} = \sigma_j^{(l)} \left(\sum_{k=0}^{p^{(l)}} \gamma_{kj}^{(l)} \beta_{kj}^{(l)} x_{ik}^{(l)} \right), j > 0 \quad (3)$$

$f(\cdot | \boldsymbol{\mu}, \phi)$ is a distribution with expectation $\boldsymbol{\mu}$ and dispersion ϕ ;

$\beta_{kj}^{(l)} \in \mathcal{R}$ are the weights (slope coefficients) for the inputs $x_{ik}^{(l)}$;

$\gamma_{kj}^{(l)} \in \{0, 1\}$ are latent indicators switching the weights on and off;

$p^{(l)}$ is the number of neurons at layer l ;

L is the total number of layers;

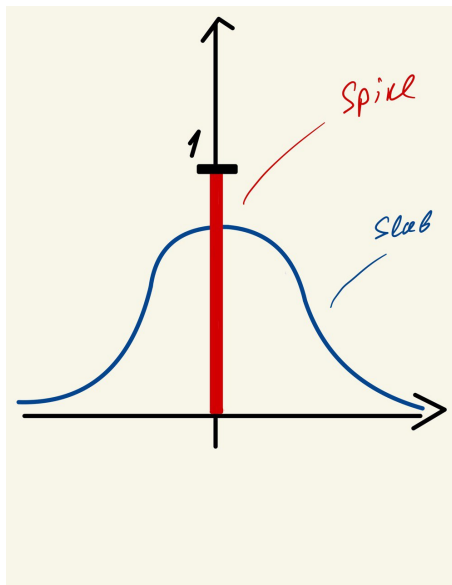
$x_{i0}^{(l+1)} = 1$ is a constant for the intercept.

Model and parameter priors

$$p(\beta_{kj}^{(l)} | \sigma_{\beta,l}^2, \gamma_{kj}^{(l)}) = \gamma_{kj}^{(l)} \mathcal{N}(0, \sigma_{\beta,l}^2) + (1 - \gamma_{kj}^{(l)}) \delta_0(\beta_{kj}^{(l)}),$$
$$p(\gamma_{kj}^{(l)}) = \text{Bernoulli}(\psi^{(l)}).$$

- $\delta_0(\cdot)$ is the delta mass or "spike" at zero;
- $\sigma_{\beta,l}^2$ is the prior variance of the weight $\beta_{kj}^{(l)}$;
- $\psi^{(l)} \in (0, 1)$ is the prior probability for including the weight $\beta_{kj}^{(l)}$.

Model and parameter priors



Model and parameter hyper priors

$$p(\sigma_{\beta,l}^{-2}) = \text{Gamma}(a_{\beta}^{(l)}, b_{\beta}^{(l)}),$$

$$p(\psi^{(l)}) = \text{Beta}(a_{\psi}^{(l)}, b_{\psi}^{(l)}).$$

- a_{β}, b_{β} hyperparameters of Gamma hyperprior for $\sigma_{\beta,l}^{-2}$;
- $a_{\psi}^{(l)}, b_{\psi}^{(l)}$ are hyperparameters of Beta hyperprior for $\psi^{(l)} \in (0, 1)$.

Inference on the model

Let:

- $\mathbf{m} = \cup_{l,j,k} \gamma_{kj}^{(l)}$ define a model itself, i.e. which weights are switched on and which are switched off;
- $\boldsymbol{\theta}|\mathbf{m} = \{\boldsymbol{\beta}, \phi|\mathbf{m}\}$, where $\boldsymbol{\beta}|\mathbf{m} = \cup_{l,j,k:\gamma_{kj}^{(l)}=1} \beta_{kj}^{(l)}$, define parameters of \mathbf{m} .

Inference on the model

Let:

- $\mathbf{m} = \cup_{l,j,k} \gamma_{kj}^{(l)}$ define a model itself, i.e. which weights are switched on and which are switched off;
- $\boldsymbol{\theta}|\mathbf{m} = \{\boldsymbol{\beta}, \phi|\mathbf{m}\}$, where $\boldsymbol{\beta}|\mathbf{m} = \cup_{l,j,k:\gamma_{kj}^{(l)}=1} \beta_{kj}^{(l)}$, define parameters of \mathbf{m} .

Goals:

- $p(\mathbf{m}, \boldsymbol{\theta}|\mathbb{D})$ posterior distribution of parameters and models;
- $p(\mathbf{m}|\mathbb{D})$ marginal posterior probabilities of the models;
- $p(\Delta|\mathbb{D})$ marginal posteriors of the parameter of interest Δ .

Inference on the model

Let:

- $\mathbf{m} = \cup_{l,j,k} \gamma_{kj}^{(l)}$ define a model itself, i.e. which weights are switched on and which are switched off;
- $\theta|\mathbf{m} = \{\beta, \phi|\mathbf{m}\}$, where $\beta|\mathbf{m} = \cup_{l,j,k:\gamma_{kj}^{(l)}=1} \beta_{kj}^{(l)}$, define parameters of \mathbf{m} .

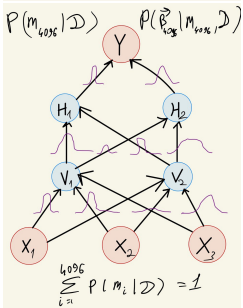
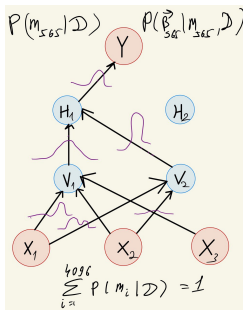
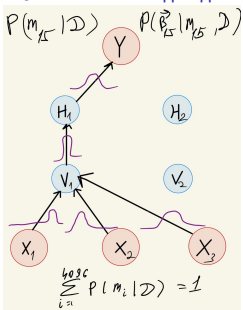
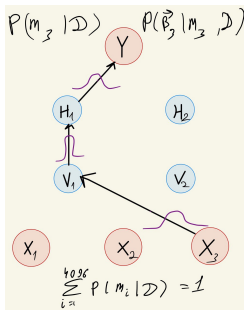
Goals:

- $p(\mathbf{m}, \theta|\mathbb{D})$ posterior distribution of parameters and models;
- $p(\mathbf{m}|\mathbb{D})$ marginal posterior probabilities of the models;
- $p(\Delta|\mathbb{D})$ marginal posteriors of the parameter of interest Δ .

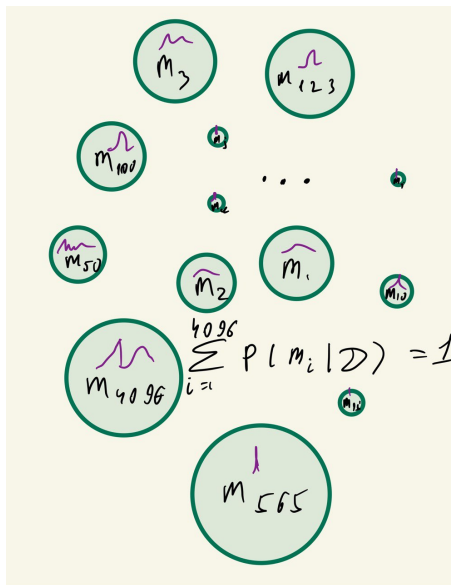
But:

- $\exists 2^q$ different models in Γ ;
- q is the number of weights in the BNN, which is huge;
- Γ is not feasible to even specify.

LBBNN. Illustrations. $q = 12 \rightarrow \|\Gamma\| = 2^q = 4096$.



The model space Γ



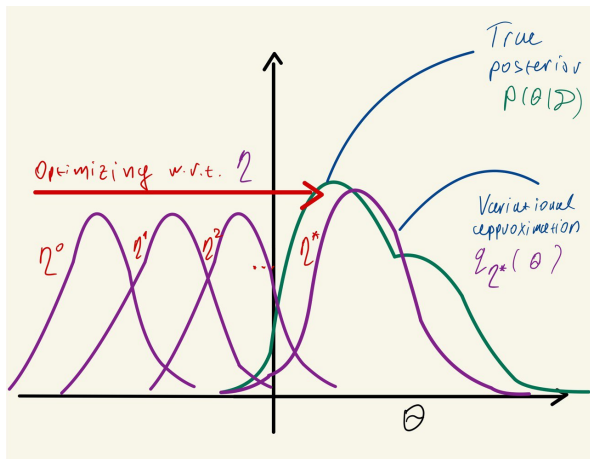
Inference possibilities

- (RJ) Markov chain Monte Carlo (exact inference) [Hubin et al., 2021];
- Laplace approximations;
- Integrated nested Laplace approximations; [Rue et al., 2009];
- **Variational inference;**
- Approximate Bayesian computation.

Variational Inference. Idea

Approximate $p(\theta|\mathbb{D})$ with $q_\eta(\theta)$ by minimizing functional divergence

$$\text{KL}(q_\eta(\theta) \| p(\theta|\mathbb{D})) = \int_{\Theta} q_\eta(\theta) \log \frac{q_\eta(\theta)}{p(\theta|\mathbb{D})} d\theta \geq 0 \text{ w.r.t. } \eta:$$



Variational inference

Posterior joint distribution $p(\boldsymbol{\theta}, \mathbf{m}|\mathbb{D})$ is approximated by combining:

- Scalable variational inference for BNN proposed by [Graves, 2011]:

$$\text{KL}(q_{\eta}(\boldsymbol{\theta}, \mathbf{m})||p(\boldsymbol{\theta}, \mathbf{m}|\mathbb{D})) = \sum_{\mathbf{m} \in \Gamma} \int_{\Theta} q_{\eta}(\boldsymbol{\theta}, \mathbf{m}) \log \frac{q_{\eta}(\boldsymbol{\theta}, \mathbf{m})}{p(\boldsymbol{\theta}, \mathbf{m}|\mathbb{D})} d\boldsymbol{\theta} \rightarrow \min_{\eta}; \quad (4)$$

Variational inference

Posterior joint distribution $p(\boldsymbol{\theta}, \mathbf{m}|\mathbb{D})$ is approximated by combining:

- Scalable variational inference for BNN proposed by [Graves, 2011]:

$$\text{KL}(q_{\eta}(\boldsymbol{\theta}, \mathbf{m})||p(\boldsymbol{\theta}, \mathbf{m}|\mathbb{D})) = \sum_{\mathbf{m} \in \Gamma} \int_{\Theta} q_{\eta}(\boldsymbol{\theta}, \mathbf{m}) \log \frac{q_{\eta}(\boldsymbol{\theta}, \mathbf{m})}{p(\boldsymbol{\theta}, \mathbf{m}|\mathbb{D})} d\boldsymbol{\theta} \rightarrow \min_{\eta} \quad (4)$$

- A mean-field variational distribution for the joint parameter-model settings for linear models introduced by [Carbonetto et al., 2012]:

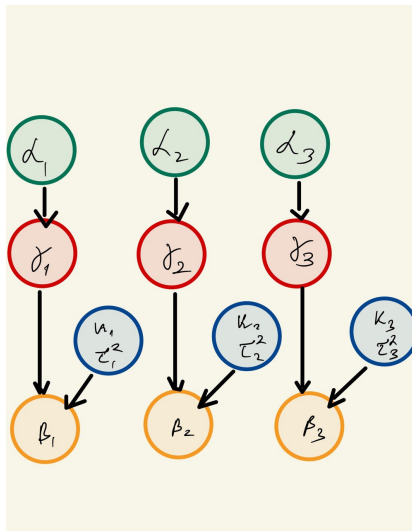
$$q_{\eta}(\boldsymbol{\theta}, \mathbf{m}) = \prod_{k,j,l} q_{\eta_{kj}^{(l)}}(\beta_{kj}^{(l)}|\gamma_{kj}^{(l)}) q_{\eta_{kj}^{(l)}}(\gamma_{kj}^{(l)}), \quad (5)$$

$$q_{\eta_{kj}^{(l)}}(\beta_{kj}^{(l)}|\gamma_{kj}^{(l)}) = \gamma_{kj}^{(l)} \mathcal{N}(\kappa_{kj}^{(l)}, \tau_{kj}^{2(l)}) + (1 - \gamma_{kj}^{(l)}) \delta_0(\beta_{kj}^{(l)}), \quad (6)$$

$$q_{\eta_{kj}^{(l)}}(\gamma_{kj}^{(l)}) = \text{Bernoulli}(\alpha_{kj}^{(l)}). \quad (7)$$

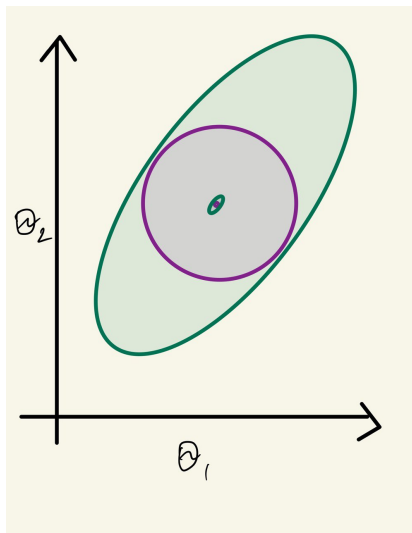
A mean-field variational distribution

A very simple approximation



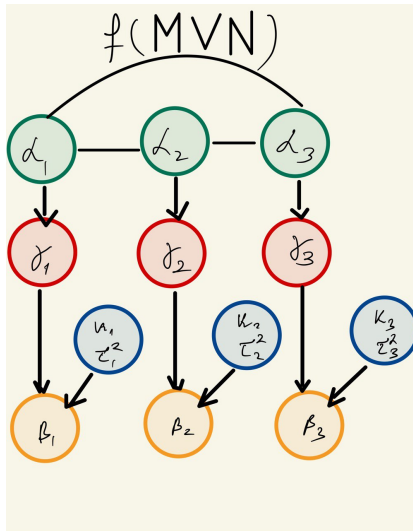
A mean-field variational distribution

But naive



Extensions of the variational distributions (MVN)

Introduce a joint Gaussian structure for transformed inclusion probabilities



Extensions of the variational distributions (MVN)

Approximation (5)-(7) assumes independence between the components, which one can argue to be unreasonable in BNNs. We proposed an extension:

$$\text{logit}(\boldsymbol{\alpha}^{(l)}) \sim \text{MVN}(\boldsymbol{\xi}^{(l)}, \boldsymbol{\Sigma}^{(l)}), \quad (8)$$

and either a full rank $\boldsymbol{\Sigma}^{(l)}$ or a low rank representation: (9)

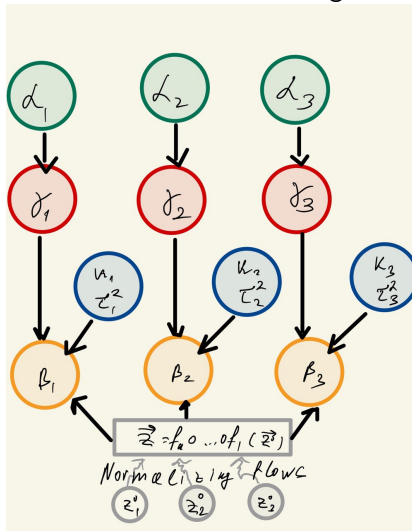
$$\boldsymbol{\Sigma}^{(l)} = \boldsymbol{F}^{(l)} \boldsymbol{F}^{(l)T} + \boldsymbol{D}^{(l)}. \quad (10)$$

For the low-ranked representation of the covariance, e.g.:

- $\boldsymbol{F}^{(l)}$ is the factor part of low-rank form of covariance matrix;
- $\boldsymbol{D}^{(l)}$ is the diagonal part of low-rank form of covariance matrix.

Extensions of the variational distributions (latent z)

Introduce a flexible structure for latent z through normalizing flows.



Extensions of the variational distributions (latent \mathbf{z})

$$q_{\boldsymbol{\eta}}(\boldsymbol{\theta}, \mathbf{m}) = \prod_{j,l} q_{\boldsymbol{\eta}_j^{(l)}}(\boldsymbol{\beta}_j^{(l)} | \gamma_j^{(l)}) q_{\boldsymbol{\eta}_j^{(l)}}(\mathbf{m}_j^{(l)}), \quad (11)$$

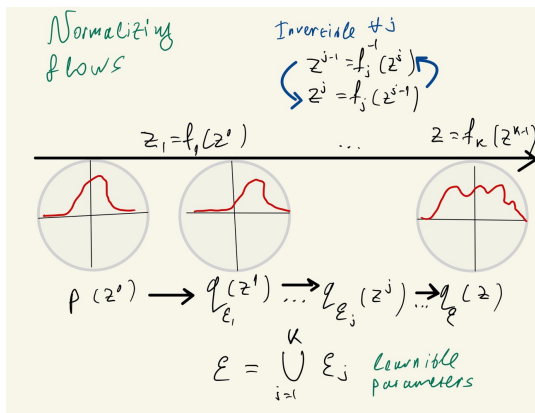
$$q_{\boldsymbol{\eta}_j^{(l)}}(\boldsymbol{\beta}_j^{(l)} | \mathbf{m}_j^{(l)}) = \int q_{\boldsymbol{\varepsilon}^{(l)}}(\mathbf{z}^{(l)}) \prod_k q_{\boldsymbol{\eta}_{kj}^{(l)}}(\boldsymbol{\beta}_{kj}^{(l)} | \gamma_{kj}^{(l)}, z_k^{(l)}) d\mathbf{z}^{(l)}, \quad (12)$$

$$q_{\boldsymbol{\eta}_{kj}^{(l)}}(\boldsymbol{\beta}_{kj}^{(l)} | \gamma_{kj}^{(l)}, z_k^{(l)}) = \gamma_{kj}^{(l)} \mathcal{N}(z_k^{(l)} \boldsymbol{\kappa}_{kj}^{(l)}, \tau_{kj}^2{}^{(l)}) + (1 - \gamma_{kj}^{(l)}) \delta_0(\boldsymbol{\beta}_{kj}^{(l)}), \quad (13)$$

$$q_{\boldsymbol{\eta}_{kj}^{(l)}}(\gamma_{kj}^{(l)}) = \text{Bernoulli}(\alpha_{kj}^{(l)}). \quad (14)$$

Autoregressive flows for latent z

$z = AF(z^0) = f_K \circ \dots \circ f_2 \circ f_1(z^0)$, is transforming a simple $z^0 \sim q_0(z^0)$ as:



Then by transformation of random variables applied K times to z^0 :

$$\log q_{\epsilon}(z_K) = \log q_0(z^0) - \sum_{k=1}^K \log \left| \det \frac{\partial f_k}{\partial z_{k-1}} \right|.$$

Technical details on $\mathbf{z} = \text{AF}(\mathbf{z}^\circ)$

$\mathbf{z} = \text{AF}(\mathbf{z}^\circ)$ from [Louizos and Welling, 2017] is transforming the standard normal \mathbf{z}° in the following way:

$$\begin{aligned}\mathbf{z}^\circ &\sim \text{MVN}(\mathbf{0}, \mathbf{I}) \\ \boldsymbol{\mu}, \mathbf{s} &= \text{NeuralNetwork}(\mathbf{z}^\circ) \\ \boldsymbol{\sigma} &= \text{sigmoid}(\mathbf{s}) \\ \mathbf{z} &= \boldsymbol{\sigma} \odot \mathbf{z}^\circ + (\mathbf{1} - \boldsymbol{\sigma}) \odot \boldsymbol{\mu},\end{aligned}\tag{15}$$

with log-determinant

$$\log \left| \frac{\partial \mathbf{z}}{\partial \mathbf{z}^\circ} \right| = \sum_{i=1}^D \log \sigma_i.\tag{16}$$

As long as the neural network in (15) is autoregressive (i.e. output dimension z_i can only depend on input dimensions up to z_{i-1}), we get a lower triangular Jacobian and therefore this simple expression for its log-determinant.

Challenges with latent \mathbf{z}

- $\int q_{\epsilon^{(l)}}(\mathbf{z}^{(l)}) \prod_k q_{\eta_{kj}^{(l)}}(\beta_{kj}^{(l)} | \gamma_{kj}^{(l)}, z_k^{(l)}) d\mathbf{z}^{(l)}$ is in general intractable;
- But $\exists q(\mathbf{z}|\boldsymbol{\theta}, \mathbf{m})$ such that $q_{\eta}(\boldsymbol{\theta}, \mathbf{m}) = \frac{q_{\eta}(\boldsymbol{\theta}, \mathbf{m}|\mathbf{z})q_{\epsilon}(\mathbf{z})}{q(\mathbf{z}|\boldsymbol{\theta}, \mathbf{m})}$;
- Then using the law of total probability

$$\text{KL}(q_{\eta}(\boldsymbol{\theta}, \mathbf{m}) || p(\boldsymbol{\theta}, \mathbf{m}|\mathbb{D})) =$$

$$\sum_{\mathbf{m} \in \mathbb{M}} \int_{\Theta} q_{\eta}(\boldsymbol{\theta}, \mathbf{m}|\mathbf{z}) q_{\epsilon}(\mathbf{z}) \log \frac{q_{\eta}(\boldsymbol{\theta}, \mathbf{m}|\mathbf{z}) q_{\epsilon}(\mathbf{z})}{p(\boldsymbol{\theta}, \mathbf{m}|\mathbb{D}) q(\mathbf{z}|\boldsymbol{\theta}, \mathbf{m})} d\mathbf{z} d\boldsymbol{\theta} =$$

$$\begin{aligned} & \mathbb{E}_{q(\boldsymbol{\theta}, \mathbf{m}, \mathbf{z})} [\text{KL}[q_{\eta}(\boldsymbol{\theta}, \mathbf{m}|\mathbf{z}) || p(\boldsymbol{\theta}, \mathbf{m})] + \log q_{\epsilon}(\mathbf{z}) - \log q(\mathbf{z}|\boldsymbol{\theta}, \mathbf{m})] \leq \\ & \mathbb{E}_{q(\boldsymbol{\theta}, \mathbf{m}, \mathbf{z})} [\text{KL}[q_{\eta}(\boldsymbol{\theta}, \mathbf{m}|\mathbf{z}) || p(\boldsymbol{\theta}, \mathbf{m})] + \log q_{\epsilon}(\mathbf{z}) - \log r_{\lambda}(\mathbf{z}|\boldsymbol{\theta}, \mathbf{m})] = \\ & \text{KL}(q(\boldsymbol{\theta}, \mathbf{m}, \mathbf{z}) || p(\boldsymbol{\theta}, \mathbf{m}|\mathbb{D}) r_{\lambda}(\mathbf{z}|\boldsymbol{\theta}, \mathbf{m})) \end{aligned}$$

- Here, we introduce use another approximation $r_{\lambda}(\mathbf{z}|\boldsymbol{\theta}, \mathbf{m})$ for $q(\mathbf{z}|\boldsymbol{\theta}, \mathbf{m})$;
- Laplace or Gaussian approximations are possible as $r_{\lambda}(\mathbf{z}|\boldsymbol{\theta}, \mathbf{m})$, but we use flexible inverse flows.

Analytic forms for $\text{KL} [q_\eta(\boldsymbol{\theta}, \mathbf{m}|\mathbf{z})||p(\boldsymbol{\theta}, \mathbf{m})]$ and $\log q_\epsilon(\mathbf{z})$

- We can show that

$$\text{KL} [q(\boldsymbol{\theta}, \mathbf{m}|\mathbf{z})||p(\boldsymbol{\theta}, \mathbf{m})] = \sum_{kj} \alpha_{q_{kj}} \left(\log \frac{\sigma_{p_{kj}}}{\sigma_{q_{kj}}} + \log \frac{\alpha_{q_{kj}}}{\alpha_{p_{kj}}} - \frac{1}{2} + \frac{\sigma_{q_{kj}}^2 + (\mu_{q_{kj}} z_{K_i} - \mu_{p_{kj}})^2}{2\sigma_{p_{kj}}^2} \right) + (1 - \alpha_{q_{kj}}) \log \frac{1 - \alpha_{q_{kj}}}{1 - \alpha_{p_{kj}}},$$

- We already saw that

$$\log q_\epsilon(\mathbf{z}) = \log q_o(\mathbf{z}^o) - \sum_{k=1}^K \log \left| \det \frac{\partial \mathbf{f}_k}{\partial \mathbf{z}_{k-1}} \right|.$$

Inverse normalizing flows for $q(\mathbf{z}|\boldsymbol{\theta}, \mathbf{m})$

- Assume $\mathbf{z}^a = \text{NF}(\mathbf{z})$ and thus $\mathbf{z} = \text{INF}(\mathbf{z}^a)$;
- Assume $r_{\lambda}(\mathbf{z}^a|\boldsymbol{\theta}, \mathbf{m}) = \prod_{i=1}^D \mathcal{N}(\tilde{\mu}_i, \tilde{\sigma}_i^2)$;
- Then by transformation formula

$$\log r(\mathbf{z}^a|\boldsymbol{\theta}, \mathbf{m}) = \log r(\mathbf{z}|\boldsymbol{\theta}, \mathbf{m}) - \sum_{t=K+1}^B \log \left| \det \frac{\partial \mathbf{g}_t}{\partial \mathbf{z}_{t-1}} \right|,$$

we know the density of $r(\mathbf{z}^a|\boldsymbol{\theta}, \mathbf{m})$ hence

$$\log r(\mathbf{z}|\boldsymbol{\theta}, \Gamma) = \log r(\mathbf{z}^a|\boldsymbol{\theta}, \Gamma) + \sum_{t=K+1}^B \log \left| \det \frac{\partial \mathbf{g}_t}{\partial \mathbf{z}_{t-1}} \right|;$$

Technical details on $\mathbf{z}^a = \text{NF}(\mathbf{z})$

- Following [Louizos and Welling, 2017] use the following flows:

$$\begin{aligned}\tilde{\boldsymbol{\mu}} &= (\mathbf{d}_1 \otimes \tanh(\mathbf{e}^T(\boldsymbol{\theta} \odot \mathbf{m}))) (1 \odot D_{\text{out}}^{-1}) \\ \log \tilde{\boldsymbol{\sigma}}^2 &= (\mathbf{d}_2 \otimes \tanh(\mathbf{e}^T(\boldsymbol{\theta} \odot \mathbf{m}))) (1 \odot D_{\text{out}}^{-1}).\end{aligned}$$

- \mathbf{d}_1 , \mathbf{d}_2 and \mathbf{e} are trainable parameters with the same shape as \mathbf{z} ;
- \otimes denotes the outer product (resulting in a matrix);
- $1 \odot D_{\text{out}}^{-1}$ denotes that we compute the mean of this matrix across its rows, resulting in a vector with the same shape as \mathbf{z} .

Evidence lower bound

Proposition

Minimization of $KL(q_{\boldsymbol{\eta}}(\boldsymbol{\theta}, \mathbf{m}) || p(\boldsymbol{\theta}, \mathbf{m} | \mathbb{D}))$ and maximization of the evidence (log marginal likelihood) lower bound (ELBO) are equivalent.

$$\mathcal{L}_{VI}(\boldsymbol{\eta}) := \sum_{\mathbf{m} \in \Gamma} \int_{\Theta} q_{\boldsymbol{\eta}}(\boldsymbol{\theta}, \mathbf{m}) \log p(\mathbb{D} | \boldsymbol{\theta}, \mathbf{m}) d\boldsymbol{\theta} - KL(q_{\boldsymbol{\eta}}(\boldsymbol{\theta}, \mathbf{m}) || p(\boldsymbol{\theta}, \mathbf{m}))$$

Evidence lower bound

Proposition

Minimization of $KL(q_{\eta}(\boldsymbol{\theta}, \mathbf{m})||p(\boldsymbol{\theta}, \mathbf{m}|\mathbb{D}))$ and maximization of the evidence (log marginal likelihood) lower bound (ELBO) are equivalent.

$$\mathcal{L}_{VI}(\boldsymbol{\eta}) := \sum_{\mathbf{m} \in \Gamma} \int_{\Theta} q_{\eta}(\boldsymbol{\theta}, \mathbf{m}) \log p(\mathbb{D}|\boldsymbol{\theta}, \mathbf{m}) d\boldsymbol{\theta} - KL(q_{\eta}(\boldsymbol{\theta}, \mathbf{m})||p(\boldsymbol{\theta}, \mathbf{m}))$$

Proof.

$$\begin{aligned} KL(q_{\eta}(\boldsymbol{\theta}, \mathbf{m})||p(\boldsymbol{\theta}, \mathbf{m}|\mathbb{D})) &= \sum_{\mathbf{m} \in \Gamma} \int_{\Theta} q_{\eta}(\boldsymbol{\theta}, \mathbf{m}) \log \frac{q_{\eta}(\boldsymbol{\theta}, \mathbf{m}) p(\mathbb{D})}{p(\mathbb{D}|\boldsymbol{\theta}, \mathbf{m}) p(\boldsymbol{\theta}, \mathbf{m})} d\boldsymbol{\theta} \\ &= \log p(\mathbb{D}) + \sum_{\mathbf{m} \in \Gamma} \int_{\Theta} q_{\eta}(\boldsymbol{\theta}, \mathbf{m}) \log \frac{q_{\eta}(\boldsymbol{\theta}, \mathbf{m})}{p(\boldsymbol{\theta}, \mathbf{m})} d\boldsymbol{\theta} - \sum_{\mathbf{m} \in \Gamma} \int_{\Theta} q_{\eta}(\boldsymbol{\theta}, \mathbf{m}) \log p(\mathbb{D}|\boldsymbol{\theta}, \mathbf{m}) d\boldsymbol{\theta} \\ &= \log p(\mathbb{D}) - \mathcal{L}_{VI}(\boldsymbol{\eta}) \geq 0. \end{aligned}$$

from which the result follows. □

Relaxations of discrete \mathbf{m} with Concrete distribution

- 1 Consider Concrete relaxations for \mathbf{m} :

$$\tilde{\gamma} = \gamma_t(\nu, \delta; \alpha) = \text{sigmoid}((\text{logit}(\alpha) - \text{logit}(\nu))/\delta), \quad \nu \sim \text{Unif}[0, 1];$$

- Here, δ is a tuning parameter with a small value;
- As δ goes to 0, $\tilde{\gamma}$ reduces to a Bernoulli(α) variable;

- 2 Consider reparametrization $\beta = \beta_t(\varepsilon; \kappa, \tau) = \kappa + \tau\varepsilon, \quad \varepsilon \sim N(0, 1)$;

- 3 Consider variational approximation

$$\mathcal{L}_{VI}^\delta(\eta) := \int_{\nu} \int_{\varepsilon} q_{\nu, \varepsilon}(\nu, \varepsilon) \left[\log p(\mathbb{D} | \beta_t(\varepsilon, \kappa, \tau), \mathbf{m}_t(\nu, \alpha, \delta)) - \log \frac{q_{\eta}(\beta_t(\varepsilon, \kappa, \tau), \mathbf{m}_t(\nu, \alpha, \delta))}{p(\beta_t(\varepsilon, \kappa, \tau), \mathbf{m}_t(\nu, \alpha, \delta))} \right] d\varepsilon d\nu;$$

- 4 Then

$$\nabla_{\eta} \mathcal{L}_{VI}^\delta(\eta) := \int_{\nu} \int_{\varepsilon} q_{\nu, \varepsilon}(\nu, \varepsilon) \frac{\partial}{\partial \eta} \left[\log p(\mathbb{D} | \beta_t(\varepsilon, \kappa, \tau), \mathbf{m}_t(\nu, \alpha, \delta)) - \log \frac{q_{\eta}(\beta_t(\varepsilon, \kappa, \tau), \mathbf{m}_t(\nu, \alpha, \delta))}{p(\beta_t(\varepsilon, \kappa, \tau), \mathbf{m}_t(\nu, \alpha, \delta))} \right] d\varepsilon d\nu.$$

- 5 And $\tilde{\nabla}_{\eta} \mathcal{L}_{VI}^\delta(\eta) =$

$$\frac{1}{M} \sum_{m=1}^M \left[\frac{n}{N} \sum_{i \in S} \nabla_{\eta} \log p(\mathbf{y}_i | \mathbf{x}_i, \beta^{(m)}, \mathbf{m}^{(m)}) - \nabla_{\eta} \log \frac{q_{\eta}(\beta^{(m)}, \mathbf{m}^{(m)})}{p(\beta^{(m)}, \mathbf{m}^{(m)})} \right]. \text{ is unbiased.}$$

Gradient estimator (unbiased for relaxed γ 's)

Proposition

Assume $(\boldsymbol{\nu}^{(t)}, \boldsymbol{\eta}^{(t)}) \sim q_{\boldsymbol{\nu}, \boldsymbol{\varepsilon}}(\boldsymbol{\nu}, \boldsymbol{\eta})$ for $t \in \{1, \dots, T\}$, $\boldsymbol{\beta}^{(t)} = \boldsymbol{\beta}_t(\boldsymbol{\varepsilon}^{(t)}, \boldsymbol{\kappa}, \boldsymbol{\tau})$, $\tilde{\mathbf{m}}^{(t)} = \mathbf{m}_t(\boldsymbol{\nu}^{(t)}, \boldsymbol{\alpha}, \delta)$ and S is a random subset of $\{1, \dots, n\}$ of size N . Then for any $\delta > 0$ an unbiased estimator for the gradient of $\mathcal{L}_{VI}^\delta(\boldsymbol{\eta})$ is given by

$$\tilde{\nabla}_{\boldsymbol{\eta}} \mathcal{L}_{VI}^\delta(\boldsymbol{\eta}) = \frac{1}{T} \sum_{t=1}^T \left[\frac{n}{N} \sum_{i \in S} \nabla_{\boldsymbol{\eta}} \log p(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\beta}^{(t)}, \mathbf{m}^{(t)}) - \nabla_{\boldsymbol{\eta}} \log \frac{q_{\boldsymbol{\eta}}(\boldsymbol{\beta}^{(t)}, \mathbf{m}^{(t)})}{p(\boldsymbol{\beta}^{(t)}, \mathbf{m}^{(t)})} \right].$$

Local reparemetrization trick (LRT)

- If within every neuron we have independent spike and slab approximations of the weights;
- It is just a mixture of Gaussians;
- Hence, we can show that the mean of their linear combination is

$$\mathbb{E}(b_j^{(l)}) = \mathbb{E} \left[\sum_{k=1}^{p^{(l)}} x_{ij}^{(l)} \gamma_{kj}^{(l)} \beta_{kj}^{(l)} \right] = \sum_{j=1}^{p^{(l)}} x_{ij}^{(l)} \alpha_{kj}^{(l)} \kappa_{kj}^{(l)}$$

- And the variance is

$$\begin{aligned} \text{Var}(b_j^{(l)}) &= \text{Var} \left[\sum_{i=k}^{p^{(l)}} x_{ij}^{(l)} \gamma_{kj}^{(l)} \beta_{kj}^{(l)} \right] = \\ &\sum_{j=1}^{p^{(l)}} x_{ij}^2 \alpha_{kj}^{(l)} (\tau_{kj}^{(l)2} + (1 - \alpha_{kj}^{(l)}) \kappa_{kj}^{(l)2}); \end{aligned}$$

- As $\alpha_{kj}^{(l)}$ converge to either 0 or 1, the mixture becomes just a unimodal Gaussian enabling direct sampling.

Model averaging

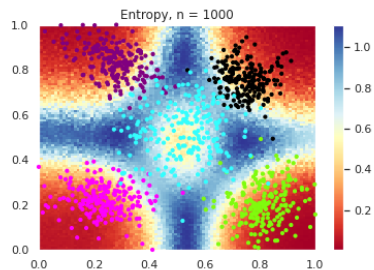
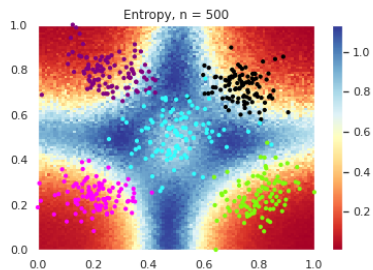
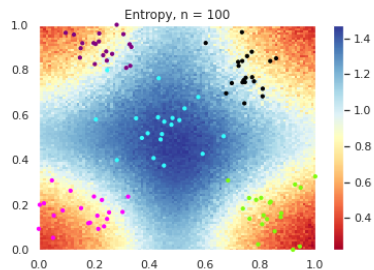
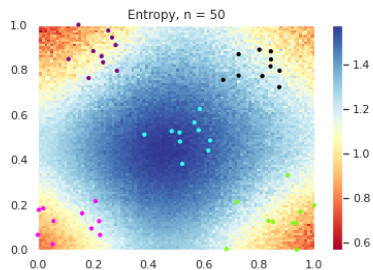
- Marginal posterior distribution of a parameter Δ (e.g. the distribution of a new observation y^* conditional on new covariates \mathbf{x}^*):

$$p(\Delta|\mathbb{D}) = \sum_{m \in \Gamma} \int_{\boldsymbol{\theta} \in \Omega_m} p(\Delta|\boldsymbol{\theta}, m, \mathbb{D}) p(\boldsymbol{\theta}, m|\mathbb{D}) d\boldsymbol{\theta}, \quad (17)$$

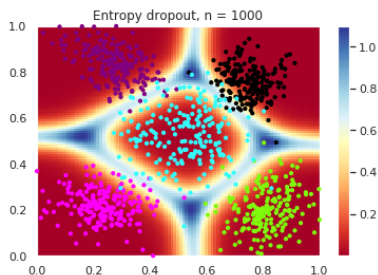
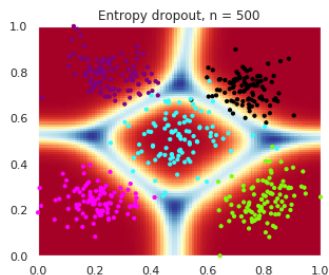
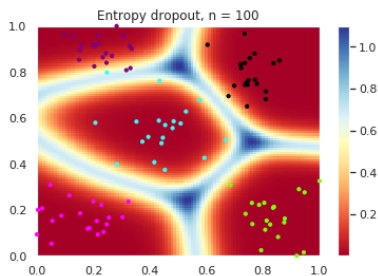
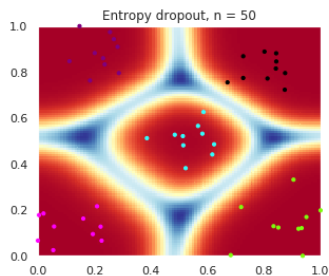
- We can now approximate it using

$$\tilde{p}(\Delta|\mathbb{D}) = \sum_{m \in \Gamma} \int_{\boldsymbol{\theta} \in \Theta_\gamma} p(\Delta|\boldsymbol{\theta}, m, \mathbb{D}) q_\eta(\boldsymbol{\theta}, m) d\boldsymbol{\theta}. \quad (18)$$

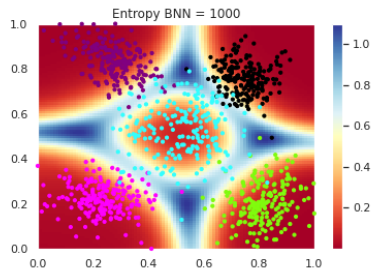
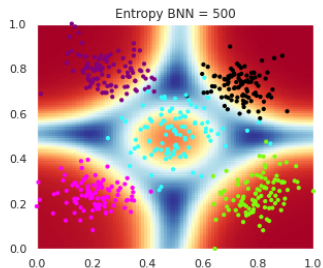
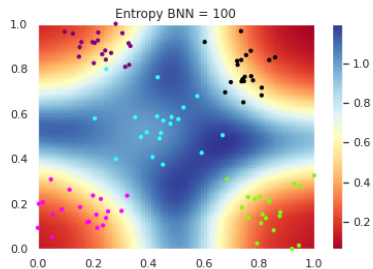
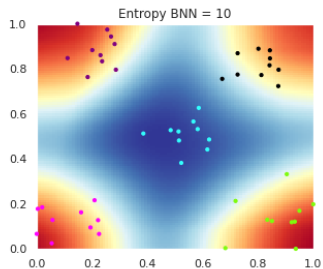
Inference with uncertainty



Do you still use dropout for uncertainty handling?



Simple dense BNNs are fine



Uncertainty FMNIST



The model for the experiments

Dense neural network with:

- ReLU activation function;
- multinomially distributed observations with 10 classes and 784 input explanatory variables (pixels);
- 2 hidden layers with 400, and 600 neurons correspondingly;
- ADAM optimizer, 250 epochs, 100 batch size.

Alternative approaches

We shall compare to:

- BNNs with a **Gaussian parameter prior**[Graves, 2011];
- BNNs with a **horseshoe prior**[Louizos et al., 2017];
- BNNs with **concrete dropout** [Gal et al., 2017].

Results. MNIST test data

method	posterior mean accuracy			model averaged accuracy			
	min	median	max	min	median	max	density
LBBNN-GP-MF	98.00	98.11	98.25	97.88	98.01	98.14	0.092
LBBNN-GP-MVN	97.60	97.80	98.00	97.60	97.80	97.90	0.180
LBBNN-GP-FLOW	98.41	98.48	98.58	98.41	98.51	98.55	0.049
BNN-GP-MF	98.20	98.40	98.50	98.20	98.30	98.50	1.000
BNN-GP-CMF	89.30	98.40	98.60	89.60	98.40	98.60	0.226
BNN-HP-MF	96.30	96.50	96.80	98.10	98.20	98.30	0.194

Table: Performance metrics on MNIST.

Results. FMNIST test data

method	posterior mean accuracy			model averaged accuracy			
	min	median	max	min	median	max	density
LBBNN-GP-MF	87.91	88.21	88.69	88.03	88.46	88.64	0.118
LBBNN-GP-MVN	86.80	87.10	87.50	87.50	87.70	87.90	0.156
LBBNN-GP-FLOW	89.57	89.74	89.93	89.48	90.76	90.07	0.046
BNN-GP-MF	88.20	88.60	88.80	89.00	89.30	89.40	1.000
BNN-GP-CMF	82.10	89.60	90.01	82.30	89.40	90.01	0.094
BNN-HP-MF	86.20	86.50	86.90	88.40	88.70	88.90	0.302

Table: Performance metrics on FMNIST

Uncertainty aware inference (0.95 threshold)

method	MNIST		FMNIST	
	decisions	accuracy	decisions	accuracy
LBBNN-GP-MF	8322	99.99	4946	99.50
LBBNN-GP-MVN	7818	100.0	4503	99.50
BNN-GP-MF	8477	99.99	5089	99.70
BNN-GP-CMF	9581	99.50	8825	94.20
BNN-HP-MF	3	100.00	181	100.0

Table: Uncertainty aware performance

Convolutional architecture (LeNet-5 style)

CNN	accuracy		density
	posterior mean	model averaged	
LBBNN-GP-MF	98.88	98.87	0.207
LBBNN-GP-FLOW	99.37	99.31	0.051

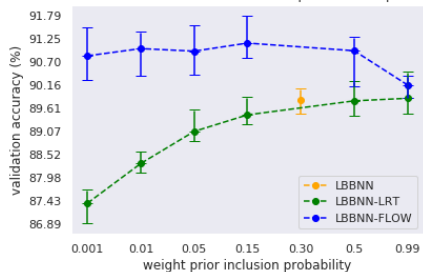
Table: CNN performance metrics on MNIST

CNN	accuracy		density
	posterior mean	model averaged	
LBBNN-GP-MF	88.30	88.24	0.209
LBBNN-GP-FLOW	90.99	91.31	0.051

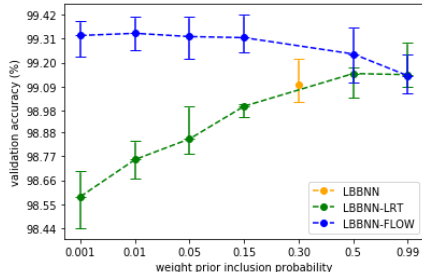
Table: CNN performance metrics on FMNIST

CNN different priors

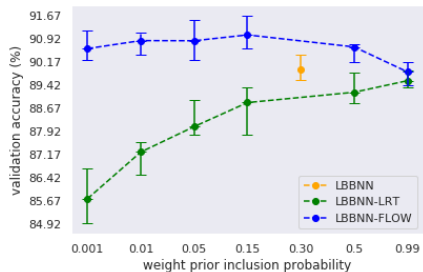
CNN architecture on FMNIST, 10 posterior samples



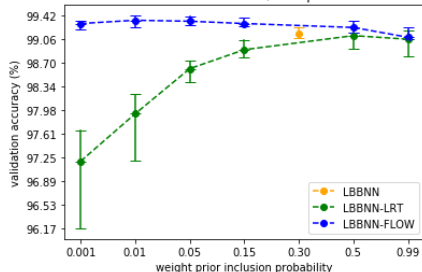
CNN architecture on MNIST, ensemble



CNN architecture on FMNIST, Posterior mean



CNN architecture on MNIST, with posterior mean



Simulation study from [Hubin and Storvik, 2018]

The response variable, Y , is generated as a logit transformation of the linear predictor in the following way,

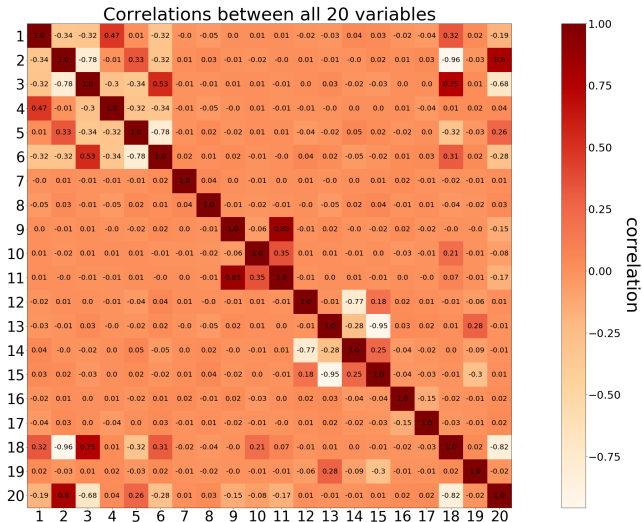
$$\eta \sim \mathcal{N}(\mathbf{X}^T \boldsymbol{\beta}, 0.5)$$

$$Y \sim \text{Bernoulli}\left(\frac{\exp(\eta)}{1 + \exp(\eta)}\right)$$

with $n = 2000$ and

$$\boldsymbol{\beta} = (-4, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1.2, 0, 37.1, 0, 0, 50, -0.00005, 10, 3, 0).$$

Correlation structure



Simulation study cont.

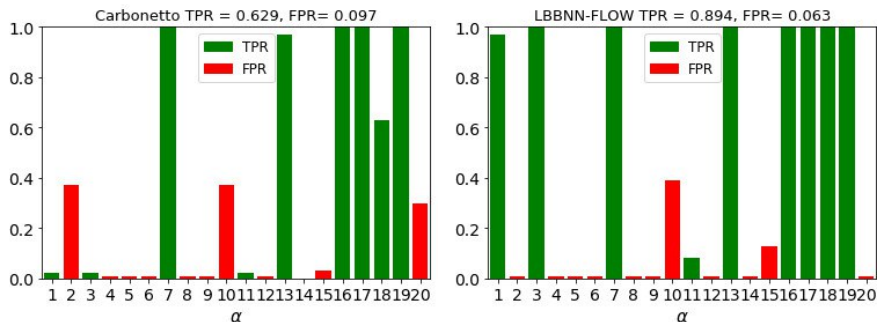





Figure: Bar plots showing true positive rate and false positive rate for the 20 different covariates on the **VARBVS**, [Carbonetto et al., 2012] (left) and **LBBNN-GP-FLOW** (right). The bars where the true weights are non-zero are colored green, and the bars where the true weights are zero are colored red.




Concluding remarks

- We develop scalable joint model-parameter approximate inference approaches in the class of BNNs;
- It allows to perform Bayesian model selection and model averaging;
- The resulting model selection often leads to drastic sparsification of BNNs with no loss of predictive power;
- Furthermore, both model selection and model averaging within our approach allow for accurate and robust handling of predictive uncertainty;
- However, the VB approach can generally be extremely biased and the ways to reduce the bias must be studied further.

References I

-  Carbonetto, P., Stephens, M., et al. (2012). Scalable variational inference for Bayesian variable selection in regression, and its accuracy in genetic association studies. *Bayesian analysis*, 7(1):73–108.
-  Gal, Y., Hron, J., and Kendall, A. (2017). Concrete dropout. In *Advances in Neural Information Processing Systems*, pages 3581–3590.
-  Graves, A. (2011). Practical variational inference for neural networks. In *Advances in Neural Information Processing Systems*, pages 2348–2356.

References II

-  Hubin, A. and Storvik, G. (2018).
Mode jumping MCMC for Bayesian variable selection in GLMM.
Computational Statistics & Data Analysis.
-  Hubin, A., Storvik, G., and Frommlet, F. (2021).
Flexible bayesian nonlinear model configuration.
Journal of Artificial Intelligence Research, 72:901–942.
-  Louizos, C., Ullrich, K., and Welling, M. (2017).
Bayesian compression for deep learning.
In *Advances in Neural Information Processing Systems*, pages 3288–3298.

References III



Louizos, C. and Welling, M. (2017).

Multiplicative normalizing flows for variational bayesian neural networks.

In International Conference on Machine Learning, pages 2218–2227. PMLR.



Rue, H., Martino, S., and Chopin, N. (2009).

Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations.

Journal of the Royal Statistical Society, 71(2):319–392.