

Marginal Likelihood Gradient for Bayesian Neural Networks

Marcin B. Tomczak MBT27@CAM.AC.UK and **Richard E. Turner** RET26@CAM.AC.UK
University of Cambridge, Cambridge, CB2 1PZ, UK

Abstract

Bayesian learning of neural networks is attractive as it can protect against over-fitting and provide automatic methods for inferring important hyperparameters by maximizing the marginal probability of the data. However, existing approaches in this vein, such as those based on variational inference, do not perform well. In this paper, we take a different approach and directly derive a practical estimator of the gradient of the marginal log-likelihood for BNNs by combining local reparametrization of the network w.r.t. the prior distribution with the self-normalized importance sampling estimator. We show promising preliminary results on a toy example and on vectorized MNIST classification where the new method results in significantly improved performance of variational inference compared to existing approaches to tune hyperparameters.

1. Introduction

We consider a supervised learning task given a data set consisting of input-target pairs $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$. We denote the observation model for data as $p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}, \boldsymbol{\beta})$ defining the likelihood $p(\mathcal{D}|\boldsymbol{\theta}, \boldsymbol{\beta}) = \prod_{i=1}^N p(\mathbf{y}_i|\mathbf{x}_i, \boldsymbol{\theta}, \boldsymbol{\beta})$ and prior $p(\boldsymbol{\theta}|\boldsymbol{\alpha})$, where $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ denote hyperparameters. Marginal likelihood of data is defined as $p(\mathcal{D}|\boldsymbol{\alpha}, \boldsymbol{\beta}) = p(\{\mathbf{y}_i\}_{i=1}^N|\{\mathbf{x}_i\}_{i=1}^N, \boldsymbol{\alpha}, \boldsymbol{\beta})$. Given a new data point $(\mathbf{x}^*, \mathbf{y}^*)$, a fully Bayesian approach to modeling makes a prediction by marginalizing out posterior over parameters $\boldsymbol{\theta}$ and hyperparameters $\boldsymbol{\alpha}, \boldsymbol{\beta}$ (MacKay, 1992a):

$$p(\mathbf{y}^*|\mathbf{x}^*, \mathcal{D}) = \int \int p(\mathbf{y}^*|\mathbf{x}^*, \boldsymbol{\theta}, \boldsymbol{\beta})p(\boldsymbol{\theta}|\mathcal{D}, \boldsymbol{\alpha}, \boldsymbol{\beta})p(\boldsymbol{\alpha}, \boldsymbol{\beta}|\mathcal{D})d\boldsymbol{\theta}d\boldsymbol{\alpha}d\boldsymbol{\beta}. \quad (1)$$

The integrals in Equation (1) inevitably carry high computational burden for non-trivial models, requiring us to recourse to approximate inference. We subsequently discuss approximations to posterior distributions $p(\boldsymbol{\theta}|\mathcal{D}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ and $p(\boldsymbol{\alpha}, \boldsymbol{\beta}|\mathcal{D})$.

Variational inference (VI). Variational inference (Barber and Bishop, 1998; Hinton and van Camp, 1993) reformulates finding an approximation to a posterior $p(\boldsymbol{\theta}|\mathcal{D}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ into an optimization problem finding a simpler distribution $q(\boldsymbol{\theta}|\boldsymbol{\lambda})$ approximating $p(\boldsymbol{\theta}|\mathcal{D})$. This is done by minimizing the divergence $D_{KL}(q(\boldsymbol{\theta}|\boldsymbol{\lambda})||p(\boldsymbol{\theta}|\mathcal{D}, \boldsymbol{\alpha}, \boldsymbol{\beta}))$, which is equivalent to maximizing the Evidence Lower Bound (ELBO), a lower bound to $\log p(\mathcal{D}|\boldsymbol{\alpha}, \boldsymbol{\beta})$, defined as:

$$\mathcal{L}(\boldsymbol{\lambda}; \boldsymbol{\alpha}, \boldsymbol{\beta}) = \sum_{i=1}^N \mathbb{E}_{q(\boldsymbol{\theta}|\boldsymbol{\lambda})} \log p(\mathbf{y}_i|\mathbf{x}_i, \boldsymbol{\theta}, \boldsymbol{\beta}) - D_{KL}(q(\boldsymbol{\theta}|\boldsymbol{\lambda})||p(\boldsymbol{\theta}|\boldsymbol{\alpha})). \quad (2)$$

Evidence framework. The evidence framework (MacKay, 1992a,b) provides a systematic framework for adaptation of hyperparameters $\boldsymbol{\alpha}, \boldsymbol{\beta}$. Equation (1) indicates that $p(\boldsymbol{\alpha}, \boldsymbol{\beta}|\mathcal{D})$, (and the marginal likelihood $p(\mathcal{D}|\boldsymbol{\alpha}, \boldsymbol{\beta})$ if the prior over $\boldsymbol{\alpha}, \boldsymbol{\beta}$ is uniform) embodies the influence of predictions arising from hyperparameters $\boldsymbol{\alpha}, \boldsymbol{\beta}$ on marginalized test predictions

$p(\mathbf{y}^*|\mathbf{x}^*, \mathcal{D})$. If the integral in Equation (1) has no closed form, it can be approximated as $p(\mathbf{y}^*|\mathbf{x}^*, \mathcal{D}) \approx \int p(\mathbf{y}^*|\mathbf{x}^*, \boldsymbol{\theta})p(\boldsymbol{\theta}|\mathcal{D}, \boldsymbol{\alpha}^*, \boldsymbol{\beta}^*)d\boldsymbol{\theta}$ where $\boldsymbol{\alpha}^*, \boldsymbol{\beta}^* = \operatorname{argmax}_{\boldsymbol{\alpha}, \boldsymbol{\beta}} p(\boldsymbol{\alpha}, \boldsymbol{\beta}|\mathcal{D})$. Relying only on the mode $\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*$ can lead to poor approximation of $p(\mathbf{y}^*|\mathbf{x}^*, \mathcal{D})$, depending on the concentration of density $p(\boldsymbol{\alpha}, \boldsymbol{\beta}|\mathcal{D})$. Estimating marginal log-likelihood $\log p(\mathcal{D}|\boldsymbol{\alpha}, \boldsymbol{\beta})$ for non-trivial models is a well-known problem (Barber and Bishop, 1998). Poorly chosen prior distribution can result in ELBO not correlating with the predictive performance of the learned posterior (Trippe and Turner, 2018; Turner and Sahani, 2011).

Bayesian Neural Networks. We consider a neural network $\mathcal{F}_{\mathbf{W}}$ represented as composition of L fully-connected layers interleaved with activation functions $\mathcal{F}_{\mathbf{W}} = \mathcal{F}_{\mathbf{W}_L} \circ \phi \circ \mathcal{F}_{\mathbf{W}_{L-1}} \dots \phi \circ \mathcal{F}_{\mathbf{W}_1}$ where ϕ denotes an activation function. We denote a forward pass with input \mathbf{x}_l through l -th layer as $\mathcal{F}_{\mathbf{W}_l}(\mathbf{x}_l) = \mathbf{W}_l \mathbf{x}_l$ where biases are handled by expanding the dimension of \mathbf{x}_l . We restrict the attention to Gaussian mean-field priors $p(\mathbf{W}_l|\boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2)$ and mean-field approximate posteriors $q(\mathbf{W}_l|\boldsymbol{\mu}_l, \boldsymbol{\sigma}_l^2)$, see Foong et al. (2019) for discussion of limitations, independent over layers and use the notation $\mathbf{W}_l \sim \mathcal{N}(\boldsymbol{\mu}_l, \boldsymbol{\sigma}_l^2)$ meaning $w_{l,ij} \sim \mathcal{N}(\mu_{l,ij}, \sigma_{l,ij}^2)$. We adopt the notation $\mathbf{W} = (\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_L)$ for weights, $\boldsymbol{\mu} = (\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_L)$, $\boldsymbol{\sigma} = (\boldsymbol{\sigma}_1^2, \boldsymbol{\sigma}_2^2, \dots, \boldsymbol{\sigma}_L^2)$ for variational parameters and $\boldsymbol{\mu}_\alpha = (\boldsymbol{\mu}_{\alpha_1}, \boldsymbol{\mu}_{\alpha_2}, \dots, \boldsymbol{\mu}_{\alpha_L})$, $\boldsymbol{\sigma}_\alpha^2 = (\boldsymbol{\sigma}_{\alpha_1}^2, \boldsymbol{\sigma}_{\alpha_1}^2, \dots, \boldsymbol{\sigma}_{\alpha_L}^2)$ for prior hyperparameters. Hyperparameter $\boldsymbol{\beta}$ belongs to an observation model, for instance, the scale of Gaussian noise $\mathcal{N}(\mathbf{y}|\mathbf{0}, \operatorname{diag}[\boldsymbol{\beta}^2])$.

Local Reparametrization Trick (LRT). The forward pass $\mathcal{F}_{\mathbf{W}_l}(\mathbf{x}_l)$ through a fully-connected layer with mean-field Gaussian distribution over weights $\mathbf{W}_l \sim \mathcal{N}(\boldsymbol{\mu}_l, \boldsymbol{\sigma}_l^2)$, can be reparametrized as (Kingma et al., 2015; Tomczak et al., 2020):

$$\mathcal{F}_{\mathbf{W}_l}(\mathbf{x}_l) = \mathcal{F}_{\boldsymbol{\mu}_l}(\mathbf{x}_l) + \sqrt{\mathcal{F}_{\boldsymbol{\sigma}_l^2}(\mathbf{x}_l)} \odot \boldsymbol{\epsilon}_l, \quad \boldsymbol{\epsilon}_l \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (3)$$

where $\sqrt{\cdot}$ is applied element-wise and the Equation (3) denotes the equality between two random variables (LHS and RHS follow the same distribution). We denote the reparametrization noise as $\boldsymbol{\epsilon} = (\boldsymbol{\epsilon}_1, \boldsymbol{\epsilon}_2, \dots, \boldsymbol{\epsilon}_L)$ and reparametrization of the network $f_{\mathbf{W}}$ resulting from applying Equation (3) to every layer l by $g_{\boldsymbol{\mu}, \boldsymbol{\sigma}^2, \boldsymbol{\epsilon}}$. Then we have that $\mathbb{E}_{q(\mathbf{W}|\boldsymbol{\mu}, \boldsymbol{\sigma}^2)} h(f_{\mathbf{W}}(\mathbf{x})) = \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} h(g_{\boldsymbol{\mu}, \boldsymbol{\sigma}^2, \boldsymbol{\epsilon}}(\mathbf{x}))$ for a continuous function h .

Approaches to adjust priors for BNNs. Hyperparameters $\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*$ can be also learned by optimizing $\mathcal{L}(\boldsymbol{\lambda}; \boldsymbol{\alpha}, \boldsymbol{\beta})$ as proxy for marginal log-likelihood since $\log p(\mathcal{D}; \boldsymbol{\alpha}, \boldsymbol{\beta}) \geq \mathcal{L}(\boldsymbol{\lambda}; \boldsymbol{\alpha}, \boldsymbol{\beta})$ (MacKay, 1999; Wu et al., 2019). The local optima of $\log p(\mathcal{D}; \boldsymbol{\alpha}, \boldsymbol{\beta})$ and $\mathcal{L}(\boldsymbol{\lambda}; \boldsymbol{\alpha}, \boldsymbol{\beta})$ can significantly vary. This is manifested when a flexible prior distribution $p(\boldsymbol{\theta}|\boldsymbol{\alpha})$ learned by optimizing $\mathcal{L}(\boldsymbol{\lambda}; \boldsymbol{\alpha}, \boldsymbol{\beta})$ collapses to the approximate posterior $q(\boldsymbol{\theta}|\boldsymbol{\lambda})$, resulting in no regularization (Blundell et al., 2015). Laplace approximation assumes $\log p(\mathcal{D}, \boldsymbol{\theta}|\boldsymbol{\alpha}, \boldsymbol{\beta})$ is a quadratic w.r.t. $\boldsymbol{\theta}$ enabling approximating $\log p(\mathcal{D}|\boldsymbol{\alpha}, \boldsymbol{\beta})$ in closed form (MacKay, 1992c), which can be used to compare different hyperparameters $\boldsymbol{\alpha}, \boldsymbol{\beta}$. Hyperparameters $\boldsymbol{\alpha}, \boldsymbol{\beta}$ can be also chosen using cross-validation (Farquhar et al., 2020; Tomczak et al., 2018), but this requires increased computational effort, limiting the practicability of this approach.

2. Methods

First, we discuss the problems with the naive application of self-normalized importance sampling (SNIS) to estimate the gradient of marginal log-likelihood of data $\nabla_{\boldsymbol{\alpha}} \log p(\mathcal{D}|\boldsymbol{\alpha}, \boldsymbol{\beta})$ for BNNs. We then demonstrate how taking the advantage of local reparametrization of

the network allows us to mitigate these problems and derive a practical estimator of the gradient of $\nabla_{\boldsymbol{\alpha}} \log p(\mathcal{D}|\boldsymbol{\alpha}, \boldsymbol{\beta})$ given a BNN with mean-field Gaussian prior. Recall that the gradient $\nabla_{\boldsymbol{\alpha}} \log p(\mathcal{D}|\boldsymbol{\alpha}, \boldsymbol{\beta})$ can be derived as (Bornschein and Bengio, 2014; Tucker et al., 2019):

$$\nabla_{\boldsymbol{\alpha}} \log p(\mathcal{D}|\boldsymbol{\alpha}, \boldsymbol{\beta}) = \mathbb{E}_{p(\mathbf{W}|\mathcal{D}, \boldsymbol{\alpha}, \boldsymbol{\beta})} \nabla_{\boldsymbol{\alpha}} \log p(\mathbf{W}|\boldsymbol{\alpha}). \quad (4)$$

See Appendix A for similar expression for $\nabla_{\boldsymbol{\beta}} \log p(\mathcal{D}|\boldsymbol{\alpha}, \boldsymbol{\beta})$ and associated derivations. A naive attempt to apply SNIS to Equation (4) results in estimator:

$$\hat{I}(\nabla_{\boldsymbol{\alpha}} \log p(\mathcal{D}|\boldsymbol{\alpha}, \boldsymbol{\beta})) = \sum_{s=1}^S \frac{w(\mathbf{W}^{(s)})}{\sum_{j=1}^S w(\mathbf{W}^{(j)})} \nabla_{\boldsymbol{\alpha}} \log p(\mathbf{W}^{(s)}|\boldsymbol{\alpha}, \boldsymbol{\beta}), \quad (5)$$

where $w(\mathbf{W}^{(s)}) = p(\mathcal{D}|\mathbf{W}^{(s)}, \boldsymbol{\beta})p(\mathbf{W}^{(s)}|\boldsymbol{\alpha})/q(\mathbf{W}^{(s)})$ and weights are proposed from $\mathbf{W}^{(s)} \sim q(\mathbf{W})$. There are several challenges associated with applying the estimator given by Equation (5) making it impractical to use with BNNs. First, sampling multiple weights $\mathbf{W}^{(s)}$ is excessively costly even for networks of moderate size as this corresponds to sampling separate copies of the network $f_{\mathbf{W}}$. Second, RHS of Equation (5) has high variance being a score gradient estimator (Williams, 1992). Third, a separate objective needs to be optimized to derive proposal $q(\mathbf{W}|\boldsymbol{\alpha})$, e.g. by minimizing a divergence between $q(\mathbf{W})$ and $p(\mathbf{W}|\mathcal{D}, \boldsymbol{\alpha}, \boldsymbol{\beta})$. Lastly, calculating weights $w(\mathbf{W}^{(s)})$ requires iterating over whole data set \mathcal{D} .

We now demonstrate that reparametrizing the network using LRT leads to the estimator free of the two first described issues. More precisely, rather than operating in the space of random weights \mathbf{W}_l , we can reparametrize the network using LRT and operate in the space of preactivation reparametrization noise $\boldsymbol{\epsilon}_l \sim N(\mathbf{0}, \mathbf{I})$. This reduces the cost of sampling, as noise $\boldsymbol{\epsilon}_l$ has significantly lower dimensionality than weight matrices \mathbf{W}_l . In addition, the derived estimator is an improvement over one given by Equation (5) in a similar way as reparametrizing (Kingma and Welling, 2014; Rezende et al., 2014; Kingma et al., 2015) improves upon the score function estimator applied to the ELBO (Paisley et al., 2012; Ranganath et al., 2014; Titsias and Lázaro-Gredilla, 2014), since path derivative usually exhibits lower variance than score function gradient (Roeder et al., 2017). We subsequently discuss how to derive better proposal distributions, reduce the cost of computing IS weights, and derive an approximation using one optimization objective in Appendix A. We now state Lemma deriving the reparametrized gradient of marginal log-likelihood of data $\log p(\mathcal{D}|\boldsymbol{\alpha}, \boldsymbol{\beta})$.

Lemma 1 (Marginal log-likelihood gradient for BNN with MF Gaussian prior)
 Given a neural network $f_{\mathbf{W}}$ with layer-wise mean-field Gaussian prior $\mathcal{N}(\mathbf{W}_l|\boldsymbol{\mu}_{\alpha_l}, \boldsymbol{\sigma}_{\alpha_l}^2)$ where $\boldsymbol{\alpha}, \boldsymbol{\beta}$ are hyperparameters, the marginal log-likelihood gradient can be computed as:

$$\nabla_{\boldsymbol{\alpha}} \log p(\mathcal{D}|\boldsymbol{\alpha}, \boldsymbol{\beta}) = \mathbb{E}_{p(\boldsymbol{\epsilon}|\mathcal{D}, \boldsymbol{\alpha}, \boldsymbol{\beta})} \sum_{i=1}^N \nabla_{\boldsymbol{\alpha}} \log p(\mathbf{y}_i|g_{\boldsymbol{\mu}_{\alpha}, \boldsymbol{\sigma}_{\alpha}^2, \boldsymbol{\epsilon}}(\mathbf{x}), \boldsymbol{\beta}), \quad (6)$$

where $p(\boldsymbol{\epsilon}|\mathcal{D}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ is a posterior distribution over layer’s preactivations resulting from applying reparametrization in Equation (3) to layers of a network $f_{\mathbf{W}}$. Substituting $\nabla_{\boldsymbol{\alpha}}$ with $\nabla_{\boldsymbol{\beta}}$ in Equation (6) yields the gradient for hyperparameters of observation model $\nabla_{\boldsymbol{\beta}} \log p(\mathcal{D}|\boldsymbol{\alpha}, \boldsymbol{\beta})$.

Note that Equation (6) still requires sampling from posterior distribution $p(\boldsymbol{\epsilon}|\mathcal{D}, \boldsymbol{\alpha}, \boldsymbol{\beta})$, where $\boldsymbol{\epsilon} = (\boldsymbol{\epsilon}_1, \boldsymbol{\epsilon}_2, \dots, \boldsymbol{\epsilon}_l)$ denotes the local reparametrization noise for consecutive layers $\mathcal{F}_{\mathbf{W}_l}$. While analyzing the posterior distribution over reparametrization noise might

be uncommon, it is well-defined: for fixed α, β and for every ϵ we define the likelihood of data $\prod_{i=1}^N p(\mathbf{y}_i | g_{\mu_\alpha, \sigma_\alpha^2, \epsilon}(\mathbf{x}), \beta)$ and hence we can calculate posterior distribution over reparametrization noise $p(\epsilon | \mathcal{D}, \alpha, \beta)$ given prior $p(\epsilon)$. The gradient in the RHS of Equation (6) resembles gradients derived for generative models (Bornschein and Bengio, 2014; Burda et al., 2015; Tucker et al., 2019). However, it is important to note not all derivations for generative models transfer to BNNs if local reparametrization is desired, as LRT requires the expression under the expectation to be only a function of layer’s outputs $\{\mathcal{F}(\mathbf{x}_l)\}_l$. For instance, the expectation $\mathbb{E}_{q(\mathbf{W})} p(\mathbf{y} | \mathbf{x}, f_{\mathbf{W}}) p(\mathbf{W}) / q(\mathbf{W})$ cannot be locally reparametrized due to the dependency on $p(\mathbf{W}) / q(\mathbf{W})$. Applying SNIS to Equation (6) by proposing from distribution $q(\epsilon)$ gives us the following estimator of $\nabla_\alpha \log p(\mathcal{D} | \mu_\alpha, \sigma_\alpha^2)$:

$$\hat{I}(\nabla_\alpha \log p(\mathcal{D} | \mu_\alpha, \sigma_\alpha^2)) = \sum_{i=1}^N \sum_{s=1}^S \frac{w(\epsilon^{(i,s)})}{\sum_{j=1}^S w(\epsilon^{(i,j)})} \nabla_\alpha \log p(\mathbf{y}_i | g_{\mu_\alpha, \sigma_\alpha^2, \epsilon^{(i,s)}}(\mathbf{x}_i)), \quad \epsilon^{(i,s)} \sim q(\epsilon), \quad (7)$$

where $w(\epsilon^{(i,s)}) = \prod_{i=1}^N p(\mathbf{y}_i | g_{\mu_\alpha, \sigma_\alpha^2, \epsilon^{(i,s)}}(\mathbf{x}), \beta) p(\epsilon^{(i,s)}) / q(\epsilon^{(i,s)})$ are normalized IS weights. We will employ the estimator of $\nabla_\alpha \log p(\mathcal{D} | \mu_\alpha, \sigma_\alpha^2)$ in Equation (7) to optimize marginal log-likelihood of data $\log p(\mathcal{D} | \mu_\alpha, \sigma_\alpha^2)$ w.r.t. prior means μ_α and variances σ_α^2 .

Better proposal distributions. We now focus on deriving better proposal distributions $q(\epsilon)$ for estimator in Equation (7). This requires finding a distribution over noise ϵ closer to posterior distribution $p(\epsilon | \mathcal{D}, \alpha, \beta)$. We define means and variances for Gaussian distributions over layer’s preactivations as $\mu_\epsilon = (\mu_{1,\epsilon}, \mu_{2,\epsilon}, \dots, \mu_{L,\epsilon})$ and $\sigma_\epsilon^2 = (\sigma_{1,\epsilon}^2, \sigma_{2,\epsilon}^2, \dots, \sigma_{L,\epsilon}^2)$. We assume factors $q(\epsilon_l | \mu_{l,\epsilon}, \sigma_{l,\epsilon}^2)$ to be mean-field Gaussian distributions and construct $q(\epsilon | \mu_\epsilon, \sigma_\epsilon^2)$ by assuming independence between layers $q(\epsilon | \mu_\epsilon, \sigma_\epsilon^2) = \prod_l q(\epsilon_l | \mu_{\epsilon_l}, \sigma_{\epsilon_l}^2)$. We learn the approximate posterior distribution $q(\epsilon | \mu_\epsilon, \sigma_\epsilon^2)$ by minimizing $D_{KL}(q(\epsilon | \mu_\epsilon, \sigma_\epsilon^2) || p(\epsilon | \mathcal{D}, \mu_\alpha, \sigma_\alpha^2))$ leading to ELBO-like expression:

$$\mathcal{L}(\mu_\epsilon, \sigma_\epsilon^2) = \sum_{i=1}^N \mathbb{E}_{q(\epsilon | \mu_\epsilon, \sigma_\epsilon^2)} \log p(\mathbf{y}_i | g_{\mu_\alpha, \sigma_\alpha^2, \epsilon}(\mathbf{x}_i)) - \sum_{l=1}^L D_{KL}(q(\epsilon_l | \mu_{\epsilon_l}, \sigma_{\epsilon_l}^2) || \mathcal{N}(\epsilon_l | \mathbf{0}, \mathbf{I})). \quad (8)$$

The optimal posterior distribution $q^*(\epsilon | \mu_\epsilon, \sigma_\epsilon^2)$ derived by optimizing the $\mathcal{L}(\mu_\epsilon, \sigma_\epsilon^2)$ depends on the prior hyperparameters $\mu_\alpha, \sigma_\alpha^2$, so in practice we will optimize $\mathcal{L}(\mu_\epsilon, \sigma_\epsilon^2)$ by performing one (or possibly a few) updates while we simultaneously learning $\mu_\alpha, \sigma_\alpha^2$. This setup of learning jointly learning model parameters $\mu_\alpha, \sigma_\alpha^2$ and auxiliary parameters $\mu_\epsilon, \sigma_\epsilon^2$ resembles algorithms developed for generative models (Kingma and Welling, 2014; Bornschein and Bengio, 2014; Burda et al., 2015; Rezende et al., 2014).

Computing weights $w(\epsilon)$. The cost of calculating weights $w(\epsilon)$ grows linearly with the size of data set \mathcal{D} . To ease the computational cost of iterating over whole data set \mathcal{D} to calculate $w(\epsilon)$, we propose to employ the “average” likelihood calculated based on subsample set of data points, similarly as Li and Turner (2016). While the theoretical properties of this approximation are not well-studied, it appears to be a sensible heuristic.

Preventing over-fitting. Adjusting hyperparameters by optimizing $\log p(\mathcal{D} | \alpha, \beta)$ (type II MLE) can be prone to over-fitting. The standard approach to mitigate over-fitting hyperparameters is to introduce hyperpriors $p(\alpha)$ and $p(\beta)$. In addition, for BNNs, the variances of units σ_l^2 can be shared across the layer l , so that $\sigma_l^2 = \sigma_l^2 \mathbf{1}$. The gradients of log hyperprior $\nabla_\alpha \log p(\alpha)$ and $\nabla_\beta \log p(\beta)$ are then added to the Equation (7).

3. Experimental demonstration

We will now empirically investigate the differences between learning hyperparameters α, β via optimizing $\log p(\mathcal{D}|\alpha, \beta)$ and $\mathcal{L}(\mu, \sigma^2; \mu_\alpha, \sigma_\alpha^2)$. We demonstrate that relying on fixed hyperparameters and optimizing ELBO w.r.t. hyperparameters α, β leads to underconfident and overconfident predictions, respectively. Employing the derived estimator of $\nabla_\alpha \log p(\mathcal{D}|\mu_\alpha, \sigma_\alpha^2)$ in Equation (7) allows us to mitigate this behavior and leads to α yielding better generalization. In our experiments we assume layer-wise Gaussian hyperprior over shared (within every layer) prior log variances $p(\log \sigma_{\alpha_l}^2) = \mathcal{N}(\log \sigma_{\alpha_l}^2 | \log[50/(N_{in} + 1)], \gamma_l^2)$, where N_{in} is the size of the layer’s input and we set $\gamma_l^2 = 500$.

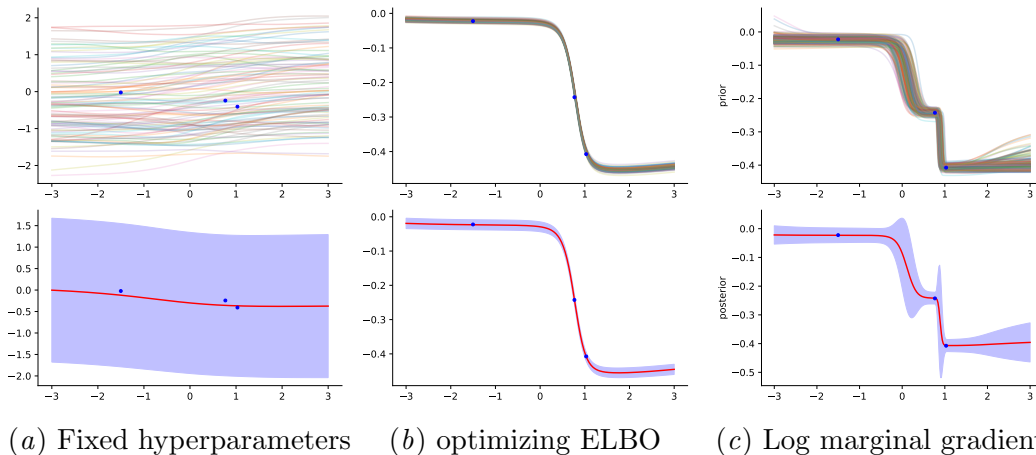


Figure 1: Comparison of different approaches to infer prior means, variances and observation noise for regression with one hidden layer BNN. Prior distribution over function (top) and posterior learned with VI (bottom).

One dimensional regression. We first consider one dimensional regression where we attempt to learn function $f(x) = (x + 1) \sin(2x + 2) + 0.3\epsilon$, $\epsilon \sim \mathcal{N}(0, 1)$ after observing 3 data points. We learn a variational posterior for two hidden layer network with 128 hidden units, tanh activations and assume Gaussian observation model $\log p(y|f_{\mathbf{W}}(x), \beta) \propto -(y - f_{\mathbf{W}}(x))^2/2\beta^2 - \log \beta$. We use ADAM optimizer (Kingma and Ba, 2014) to optimize $\mathcal{L}(\mu, \sigma^2; \mu_\alpha, \sigma_\alpha^2)$. We consider three approaches: (a) use fixed value of observation noise $\beta = 1$, zero prior unit means μ and fixed variances $\sigma_\alpha^2 = 1/N_{in}$ (Neal, 1996), (b) optimize $\mathcal{L}(\mu, \sigma^2; \mu_\alpha, \sigma_\alpha^2, \beta)$ w.r.t. both variational parameters μ, σ^2 , and prior hyperparameters $\mu_\alpha, \sigma_\alpha^2, \beta$, and (c) learn $\mu_\alpha, \sigma_\alpha^2, \beta$ using marginal log-likelihood gradient given by Equation (7) and then learn variational posterior $q(\mu, \sigma^2)$ by optimizing $\mathcal{L}(\mu, \sigma^2; \mu_\alpha, \sigma_\alpha^2, \beta)$. We report prior and posterior predictive distribution in Figure 1. First, in Figure 1(a) we observe that MF-VI with predefined variances of units can severely underfit data if prior variances of the units are not tuned. This happens when the optimization of $\mathcal{L}(\mu, \sigma^2; \mu_\alpha, \sigma_\alpha^2, \beta)$ is driven by reducing complexity penalty $D_{KL}(q(\theta|\lambda)||p(\theta|\alpha))$ and variational posterior explains data \mathcal{D} by noise. Next, in Figure 1(b) we see that learning prior hyperparameters by optimizing $\mathcal{L}(\mu, \sigma^2; \mu_\alpha, \sigma_\alpha^2, \beta)$ leads to an improved performance, but results in overconfident predictive distribution. This happens as optimizing prior $p(\mathbf{W}|\mu_\alpha, \sigma_\alpha^2)$ using Equation (2)

can set it very close to posterior $q(\mathbf{W}|\boldsymbol{\mu}, \boldsymbol{\sigma}^2)$, yielding almost no regularization. Lastly, we see from Figure 1(c) that using the derived marginal log-likelihood gradient results in a prior explaining the data \mathcal{D} with subsequent application of VI yielding reasonable performance.

Vectorized MNIST classification. Next we investigate if adjusting hyperparameters leads to improved performance of Gaussian MF-VI on a vectorized MNIST classification task, a standard benchmark for Bayesian deep learning (Blundell et al., 2015). We learn a two hidden layer network with 400 hidden units and ReLU activations (Nair and Hinton, 2010) and report learning curves as optimization proceeds in Figure 2. We use the same hyperpriors for all approaches and compare using predefined hyperparameters with learning them via optimizing ELBO and using the derived marginal log-likelihood gradient.

Overall, we observe similar patterns as in 1D regression experiment. Using predefined set of hyperparameters yields underfitting with test NLL 0.095 and test error rate 2.49%, as optimizing ELBO constrains $q(\mathbf{W}|\boldsymbol{\mu}, \boldsymbol{\sigma}^2)$ to prior distribution incapable of modeling data. Tuning prior by optimizing ELBO leads to over-fitting (test NLL 0.177 and test error rate 1.57%), as prior distribution $p(\mathbf{W}|\boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2)$ collapse to posterior $q(\mathbf{W}|\boldsymbol{\mu}, \boldsymbol{\sigma}^2)$ leading to insufficient regularization. Using marginal log-likelihood gradient given by Equation (7) optimizes prior $p(\mathbf{W}|\boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2)$ directly to fit the data and introduces sufficient regularization when optimizing ELBO, leading to the best test log likelihood 0.047 and test error rate 1.37%.

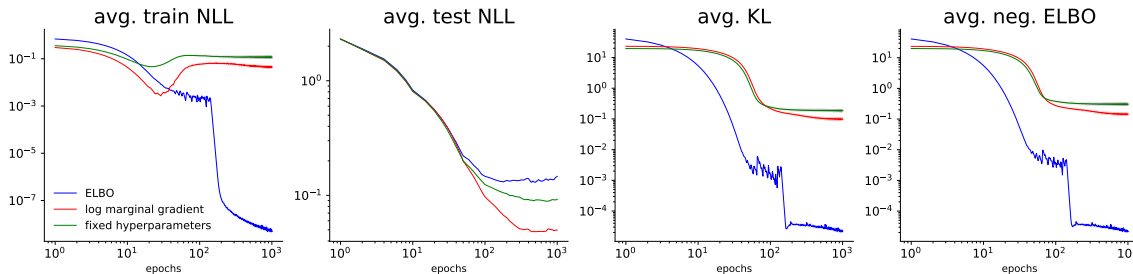


Figure 2: Different approaches to infer prior means and shared variances across layers for MF-VI on vectorized MNIST classification. Using marginal log-likelihood gradient given by Equation (7) (red curve) leads to the best predictive performance.

4. Conclusions

We have derived a pragmatic algorithm adjusting hyperparameters of BNNs with mean-field Gaussian priors by optimizing marginal log-likelihood of data. To this end, we applied local reparametrization of the network combined with self-normalized importance sampling. The optimized priors yielded predictive distributions explaining the data and subsequent application of VI resulted in posterior distributions with good predictive performance. We hope our work will catalyze research on benefits of reparametrization and algorithms inferring hyperparameters for BNNs.

References

- D. Barber and Christopher Bishop. Ensemble learning in bayesian neural networks. In *Generalization in Neural Networks and Machine Learning*, pages 215–237. Springer Verlag, January 1998. URL <https://www.microsoft.com/en-us/research/publication/ensemble-learning-in-bayesian-neural-networks/>.
- Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural networks. In *Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37*, ICML’15, page 1613–1622. JMLR.org, 2015.
- Jörg Bornschein and Yoshua Bengio. Reweighted wake-sleep. 2014.
- Yuri Burda, Roger Grosse, and Ruslan Salakhutdinov. Importance weighted autoencoders. 2015.
- Sebastian Farquhar, Michael A. Osborne, and Yarin Gal. Radial bayesian neural networks: Beyond discrete support in large-scale bayesian deep learning. In Silvia Chiappa and Roberto Calandra, editors, *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 1352–1362. PMLR, 26–28 Aug 2020.
- Andrew YK Foong, David R Burt, Yingzhen Li, and Richard E Turner. On the expressiveness of approximate inference in bayesian neural networks. *arXiv e-prints*, pages arXiv–1909, 2019.
- Geoffrey E. Hinton and Drew van Camp. Keeping the neural networks simple by minimizing the description length of the weights. In *Proceedings of the Sixth Annual Conference on Computational Learning Theory*, COLT ’93, page 5–13, New York, NY, USA, 1993. Association for Computing Machinery. ISBN 0897916115. doi: 10.1145/168304.168306. URL <https://doi.org/10.1145/168304.168306>.
- Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *International Conference on Learning Representations*, 12 2014.
- Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014. URL <http://arxiv.org/abs/1312.6114>.
- Diederik P. Kingma, Tim Salimans, and Max Welling. Variational dropout and the local reparameterization trick. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2*, NIPS’15, page 2575–2583, Cambridge, MA, USA, 2015. MIT Press.
- Yingzhen Li and Richard E Turner. Rényi divergence variational inference. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 1073–1081. Curran Associates, Inc., 2016. URL <http://papers.nips.cc/paper/6208-renyi-divergence-variational-inference.pdf>.

- David J. C. MacKay. Bayesian interpolation. *Neural Comput.*, 4(3):415–447, May 1992a. ISSN 0899-7667. doi: 10.1162/neco.1992.4.3.415. URL <https://doi.org/10.1162/neco.1992.4.3.415>.
- David J. C. MacKay. The evidence framework applied to classification networks. 4(5): 720–736, September 1992b. ISSN 0899-7667. doi: 10.1162/neco.1992.4.5.720. URL <https://doi.org/10.1162/neco.1992.4.5.720>.
- David J. C. MacKay. A practical bayesian framework for backpropagation networks. *Neural Comput.*, 4(3):448–472, May 1992c. ISSN 0899-7667. doi: 10.1162/neco.1992.4.3.448. URL <https://doi.org/10.1162/neco.1992.4.3.448>.
- David J. C. MacKay. Comparison of approximate methods for handling hyperparameters. *Neural Comput.*, 11(5):1035–1068, July 1999. ISSN 0899-7667. doi: 10.1162/089976699300016331. URL <https://doi.org/10.1162/089976699300016331>.
- Vinod Nair and Geoffrey E. Hinton. Rectified linear units improve restricted boltzmann machines. ICML’10, page 807–814, Madison, WI, USA, 2010. Omnipress. ISBN 9781605589077.
- Radford M. Neal. *Bayesian Learning for Neural Networks*. Springer-Verlag, Berlin, Heidelberg, 1996. ISBN 0387947248.
- John Paisley, David M. Blei, and Michael I. Jordan. Variational bayesian inference with stochastic search. In *Proceedings of the 29th International Conference on International Conference on Machine Learning*, ICML’12, page 1363–1370, Madison, WI, USA, 2012. Omnipress. ISBN 9781450312851.
- Rajesh Ranganath, Sean Gerrish, and David Blei. Black Box Variational Inference. In Samuel Kaski and Jukka Corander, editors, *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Statistics*, volume 33 of *Proceedings of Machine Learning Research*, pages 814–822, Reykjavik, Iceland, 22–25 Apr 2014. PMLR. URL <http://proceedings.mlr.press/v33/ranganath14.html>.
- Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In Eric P. Xing and Tony Jebara, editors, *Proceedings of the 31st International Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pages 1278–1286, Beijing, China, 22–24 Jun 2014. PMLR. URL <http://proceedings.mlr.press/v32/rezende14.html>.
- Geoffrey Roeder, Yuhuai Wu, and David K Duvenaud. Sticking the landing: Simple, lower-variance gradient estimators for variational inference. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 6925–6934. Curran Associates, Inc., 2017.
- Michalis Titsias and Miguel Lázaro-Gredilla. Doubly stochastic variational bayes for non-conjugate inference. In Eric P. Xing and Tony Jebara, editors, *Proceedings of the 31st International Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pages 1278–1286, Beijing, China, 22–24 Jun 2014. PMLR. URL <http://proceedings.mlr.press/v32/titsias14.html>.

- Learning Research*, pages 1971–1979, Beijing, China, 22–24 Jun 2014. PMLR. URL <http://proceedings.mlr.press/v32/titsias14.html>.
- Marcin B. Tomczak, Siddharth Swaroop, and Richard E. Turner. Neural network ensembles and variational inference revisited. In *1st Symposium on Advances in Approximate Bayesian Inference*, pages 1–11, 2018.
- Marcin B. Tomczak, Siddharth Swaroop, and Richard E. Turner. Efficient low rank gaussian variational inference for neural networks. In *Advances in Neural Information Processing Systems*, volume 33, 2020. URL <https://papers.nips.cc/paper/2020/file/310cc7ca5a76a446f85c1a0d641ba96d-Paper.pdf>.
- Brian Trippe and Richard Turner. Overpruning in variational bayesian neural networks. *arXiv preprint arXiv:1801.06230*, 2018.
- George Tucker, Dieterich Lawson, Shixiang Gu, and Chris J. Maddison. Doubly reparameterized gradient estimators for monte carlo objectives. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=HkG3e205K7>.
- R. E. Turner and M. Sahani. Two problems with variational expectation maximisation for time-series models. In D. Barber, T. Cemgil, and S. Chiappa, editors, *Bayesian Time series models*, chapter 5, pages 109–130. Cambridge University Press, 2011.
- Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Mach. Learn.*, 8(3–4):229–256, May 1992. ISSN 0885-6125. doi: 10.1007/BF00992696. URL <https://doi.org/10.1007/BF00992696>.
- Anqi Wu, Sebastian Nowozin, Edward Meeds, Richard E. Turner, Jose Miguel Hernandez-Lobato, and Alexander L. Gaunt. Deterministic variational inference for robust bayesian neural networks. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=B1l08oAct7>.

Appendix A.

We provide derivations of $\nabla_{\alpha} \log p(\mathcal{D}|\alpha, \beta)$ and $\nabla_{\beta} \log p(\mathcal{D}|\alpha, \beta)$ in Equation (6).

$$\nabla_{\alpha} \log p(\mathcal{D}|\alpha, \beta) = \frac{1}{p(\mathcal{D}|\alpha, \beta)} \nabla_{\alpha} p(\mathcal{D}|\alpha, \beta) \quad (9)$$

$$= \frac{1}{p(\mathcal{D}|\alpha, \beta)} \nabla_{\alpha} \int p(\mathcal{D}|\mathbf{W}, \beta) p(\mathbf{W}|\alpha) d\mathbf{W} \quad (10)$$

$$= \frac{1}{\log p(\mathcal{D}|\alpha, \beta)} \int \nabla_{\alpha} p(\mathbf{W}|\alpha) p(\mathcal{D}|\mathbf{W}, \beta) d\mathbf{W} \quad (11)$$

$$= \frac{1}{p(\mathcal{D}|\alpha, \beta)} \int p(\mathbf{W}|\alpha) \nabla_{\alpha} \log p(\mathbf{W}|\alpha) p(\mathcal{D}|\mathbf{W}, \beta) d\mathbf{W} \quad (12)$$

$$= \int \frac{p(\mathcal{D}|\mathbf{W}, \beta) p(\mathbf{W}|\alpha)}{p(\mathcal{D}|\alpha, \beta)} \nabla_{\alpha} \log p(\mathbf{W}|\alpha) d\mathbf{W} \quad (13)$$

$$= \mathbb{E}_{p(\mathbf{W}|\mathcal{D}, \alpha, \beta)} \nabla_{\alpha} \log p(\mathbf{W}|\alpha). \quad (14)$$

$$\nabla_{\beta} \log p(\mathcal{D}|\alpha, \beta) = \frac{1}{p(\mathcal{D}|\alpha, \beta)} \nabla_{\beta} p(\mathcal{D}|\alpha, \beta) \quad (15)$$

$$= \frac{1}{p(\mathcal{D}|\alpha, \beta)} \nabla_{\beta} \int p(\mathcal{D}|\mathbf{W}, \beta) p(\mathbf{W}|\alpha) d\mathbf{W} \quad (16)$$

$$= \frac{1}{p(\mathcal{D}|\alpha, \beta)} \int p(\mathbf{W}|\alpha) \nabla_{\beta} p(\mathcal{D}|\mathbf{W}, \beta) d\mathbf{W} \quad (17)$$

$$= \frac{1}{p(\mathcal{D}|\alpha, \beta)} \int p(\mathbf{W}|\alpha) p(\mathcal{D}|\mathbf{W}, \beta) \nabla_{\beta} \log p(\mathcal{D}|\mathbf{W}, \beta) d\mathbf{W} \quad (18)$$

$$= \int \frac{p(\mathcal{D}|\mathbf{W}, \beta) p(\mathbf{W}|\alpha)}{p(\mathcal{D}|\alpha, \beta)} \nabla_{\beta} \log p(\mathcal{D}|\mathbf{W}, \beta) \quad (19)$$

$$= \mathbb{E}_{p(\mathbf{W}|\mathcal{D}, \alpha, \beta)} \nabla_{\beta} \log p(\mathcal{D}|\mathbf{W}, \beta). \quad (20)$$

We now present the proof of Lemma 1. The derivation uses the reparametrization of the network $f_{\mathbf{W}}$ into $g_{\mu, \sigma^2, \epsilon}$ together with applying log trick $\nabla f(x) = f(x) \nabla \log f(x)$.

Proof Derivation of Lemma 1.

$$\nabla_{\alpha} \log p(\mathcal{D}|\mu_{\alpha}, \sigma_{\alpha}^2) = \quad (21)$$

$$= \frac{1}{p(\mathcal{D}|\mu_{\alpha}, \sigma_{\alpha}^2)} \nabla_{\alpha} p(\mathcal{D}|\mu_{\alpha}, \sigma_{\alpha}^2) \quad (22)$$

$$= \frac{1}{p(\mathcal{D}|\mu_{\alpha}, \sigma_{\alpha}^2)} \nabla_{\alpha} \int p(\mathcal{D}|\mathbf{W}, \beta) p(\mathbf{W}|\mu_{\alpha}, \sigma_{\alpha}^2) d\mathbf{W} \quad (23)$$

$$= \frac{1}{p(\mathcal{D}|\mu_{\alpha}, \sigma_{\alpha}^2)} \nabla_{\alpha} \int \prod_{i=1}^N p(\mathbf{y}_i | f_{\mathbf{W}}(\mathbf{x}_i), \beta) p(\mathbf{W}|\mu_{\alpha}, \sigma_{\alpha}^2) d\mathbf{W} \quad (24)$$

$$= \frac{1}{p(\mathcal{D}|\mu_{\alpha}, \sigma_{\alpha}^2)} \nabla_{\alpha} \int \prod_{i=1}^N p(\mathbf{y}_i | g_{\mu_{\alpha}, \sigma_{\alpha}^2, \epsilon}(\mathbf{x}_i), \beta) p(\epsilon) d\epsilon \quad (25)$$

$$= \frac{1}{p(\mathcal{D}|\boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2)} \int p(\boldsymbol{\epsilon}) \nabla_\alpha \prod_{i=1}^N p(\mathbf{y}_i | g_{\boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2, \boldsymbol{\epsilon}}(\mathbf{x}_i), \boldsymbol{\beta}) d\boldsymbol{\epsilon} \quad (26)$$

$$= \frac{1}{p(\mathcal{D}|\boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2)} \int p(\boldsymbol{\epsilon}) \prod_{i=1}^N p(\mathbf{y}_i | g_{\boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2, \boldsymbol{\epsilon}}(\mathbf{x}_i), \boldsymbol{\beta}) \nabla_\alpha \log \prod_{i=1}^N p(\mathbf{y}_i | g_{\boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2, \boldsymbol{\epsilon}}(\mathbf{x}_i), \boldsymbol{\beta}) d\boldsymbol{\epsilon} \quad (27)$$

$$= \frac{1}{p(\mathcal{D}|\boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2)} \int p(\boldsymbol{\epsilon}) \prod_{i=1}^N p(\mathbf{y}_i | g_{\boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2, \boldsymbol{\epsilon}}(\mathbf{x}_i), \boldsymbol{\beta}) \nabla_\alpha \sum_{i=1}^N \log p(\mathbf{y}_i | g_{\boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2, \boldsymbol{\epsilon}}(\mathbf{x}_i), \boldsymbol{\beta}) d\boldsymbol{\epsilon} \quad (28)$$

$$= \int \frac{p(\boldsymbol{\epsilon}) \prod_{i=1}^N p(\mathbf{y}_i | g_{\boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2, \boldsymbol{\epsilon}}(\mathbf{x}_i))}{p(\mathcal{D}|\boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2)} \sum_{i=1}^N \nabla_\alpha \log p(\mathbf{y}_i | g_{\boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2, \boldsymbol{\epsilon}}(\mathbf{x}_i), \boldsymbol{\beta}) d\boldsymbol{\epsilon} \quad (29)$$

$$= \int \frac{p(\boldsymbol{\epsilon}) p(\mathcal{D}|\boldsymbol{\epsilon}, \boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2)}{p(\mathcal{D}|\boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2)} \sum_{i=1}^N \nabla_\alpha \log p(\mathbf{y}_i | g_{\boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2, \boldsymbol{\epsilon}}(\mathbf{x}_i), \boldsymbol{\beta}) d\boldsymbol{\epsilon} \quad (30)$$

$$= \mathbb{E}_{p(\boldsymbol{\epsilon}|\mathcal{D}, \boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2)} \sum_{i=1}^N \nabla_\alpha \log p(\mathbf{y}_i | g_{\boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2, \boldsymbol{\epsilon}}(\mathbf{x}_i), \boldsymbol{\beta}) \quad (31)$$

Following the same derivations for ∇_β yields

$$\nabla_\beta \log p(\mathcal{D}|\boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2) = \mathbb{E}_{p(\boldsymbol{\epsilon}|\mathcal{D}, \boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2)} \sum_{i=1}^N \nabla_\beta \log p(\mathbf{y}_i | g_{\boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2, \boldsymbol{\epsilon}}(\mathbf{x}_i), \boldsymbol{\beta}), \quad (32)$$

■

Approximate scheme using one objective. Approximating weights with uniform distribution $w(\boldsymbol{\epsilon}^{(i,s)}) = \frac{1}{S}$ can be used to formulate the objective that can be used to learn both $\boldsymbol{\mu}_\epsilon, \boldsymbol{\sigma}_\epsilon^2$ and $\boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2$:

$$\begin{aligned} \tilde{\mathcal{L}}_1(\boldsymbol{\mu}_\epsilon, \boldsymbol{\sigma}_\epsilon^2, \boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2) &\approx \\ &\sum_{i=1}^N \mathbb{E}_{\boldsymbol{\epsilon}'} \log p(\mathbf{y}_i | g_{\boldsymbol{\mu}_\alpha, \boldsymbol{\sigma}_\alpha^2, \boldsymbol{\mu}_\epsilon + \boldsymbol{\sigma}_\epsilon \boldsymbol{\epsilon}'}(\mathbf{x}_i)) - \sum_{l=1}^L D_{KL}(q(\boldsymbol{\epsilon}_l | \boldsymbol{\mu}_{\boldsymbol{\epsilon}_l}, \boldsymbol{\sigma}_{\boldsymbol{\epsilon}_l}^2) || \mathcal{N}(\boldsymbol{\epsilon}_l | \mathbf{0}, \mathbf{I})) + \log p(\boldsymbol{\alpha}, \boldsymbol{\beta}). \end{aligned} \quad (33)$$