
Bayesian Uncertainty Estimation for Batch Normalized Deep Networks

Mattias Teye^{1,2*} Hossein Azizpour^{1*} Kevin Smith^{1,3}

Abstract

We show that training a deep network using batch normalization is equivalent to approximate inference in Bayesian models. We further demonstrate that this finding allows us to make meaningful estimates of the model uncertainty using conventional architectures, without modifications to the network or the training procedure. Our approach is thoroughly validated by measuring the quality of uncertainty in a series of empirical experiments on different tasks. It outperforms baselines with strong statistical significance, and displays competitive performance with recent Bayesian approaches.

1. Introduction

Deep learning has dramatically advanced the state of the art in a number of domains. Despite their unprecedented discriminative power, deep networks are prone to make mistakes. Nevertheless, they can already be found in settings where errors carry serious repercussions such as autonomous vehicles (Chen et al., 2016) and high frequency trading. We can soon expect automated systems to screen for various types of cancer (Esteva et al., 2017; Shen, 2017) and diagnose biopsies (Djuric et al., 2017). As autonomous systems based on deep learning are increasingly deployed in settings with the potential to cause physical or economic harm, we need to develop a better understanding of when we can be confident in the estimates produced by deep networks, and when we should be less certain.

Standard deep learning techniques used for supervised learning lack methods to account for uncertainty in the model. This can be problematic when the network encounters conditions it was not exposed to during training,

or if the network is confronted with adversarial examples (Goodfellow et al., 2014). When exposed to data outside the distribution it was trained on, the network is forced to extrapolate, which can lead to unpredictable behavior.

If the network can provide information about its uncertainty in addition to its point estimate, disaster may be avoided. In this work, we focus on estimating such predictive uncertainties in deep networks (Figure 1).

The Bayesian approach provides a theoretical framework for modeling uncertainty (Ghahramani, 2015), which has prompted several attempts to extend neural networks (NN) into a Bayesian setting. Most notably, Bayesian neural networks (BNNs) have been studied since the 1990's (Neal, 2012), but do not scale well and struggle to compete with modern deep learning architectures. Recently, (Gal & Ghahramani, 2015) developed a practical solution to obtain uncertainty estimates by casting *dropout* training in conventional deep networks as a Bayesian approximation of a Gaussian Process (its correspondence to a general approximate Bayesian model was shown in (Gal, 2016)). They showed that *any network* trained with dropout is an approximate Bayesian model, and uncertainty estimates can be obtained by computing the variance on multiple predictions with different dropout masks.

The inference in this technique, called *Monte Carlo Dropout* (MCDO), has an attractive quality: it can be applied to any pre-trained networks with dropout layers. Uncertainty estimates come (nearly) for free. However, not all architectures use dropout, and most modern networks have adopted other regularization techniques. *Batch normalization* (BN), in particular, has become widespread thanks to its ability to stabilize learning with improved generalization (Ioffe & Szegedy, 2015).

An interesting aspect of BN is that the mini-batch statistics used for training each iteration depend on randomly selected batch members. We exploit this stochasticity and show that training using batch normalization, like dropout, can be cast as an approximate Bayesian inference. We demonstrate how this finding allows us to make meaningful estimates of the model uncertainty in a technique we call *Monte Carlo Batch Normalization* (MCBN) (Figure 1). The method we propose can be applied to any network using standard batch normalization.

* Co-first authorship ¹School of Electrical Engineering and Computer Science, KTH Royal Institute of Technology, Stockholm, Sweden ²Current address: Electronic Arts, SEED, Stockholm, Sweden. This work was carried out at Budbee AB. ³Science for Life Laboratory. Correspondence to: Kevin Smith <ksmith@kth.se>.

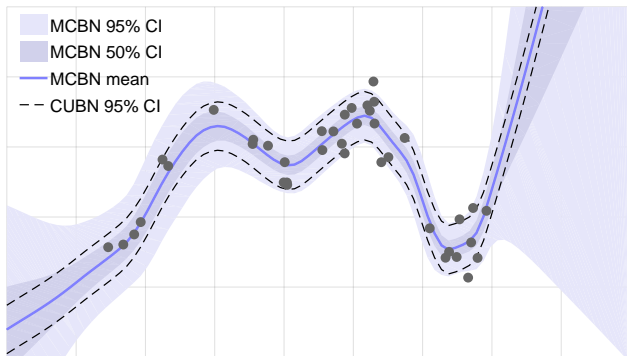


Figure 1. Training a deep network using batch normalization is equivalent to approximate inference in Bayesian models. Thus, uncertainty estimates can be obtained from any network using BN through a simple procedure. At inference, several mini-batches are constructed by taking random samples to accompany the query. The mean and variance of the outputs are used to estimate the predictive distribution (MCBN). Here, we show results on a toy dataset from a network with three hidden layers (30 units per layer). Training data is depicted as dots. The solid line is the predictive mean of 500 stochastic forward passes and the shaded areas represent the model’s uncertainty. The dashed lines depict a minimal baseline for uncertainty (CUBN), see Section 4.1.

We validate our approach by empirical experiments on a variety of datasets and tasks, including regression and image classification. We measure uncertainty quality relative to a baseline of fixed uncertainty, and show that MCBN outperforms the baseline on nearly all datasets with strong statistical significance. We also show that the uncertainty quality of MCBN is on par with other recent approximate Bayesian networks.

2. Related Work

Bayesian models provide a natural framework for modeling uncertainty, and several approaches have been developed to adapt NNs to Bayesian reasoning. A common approach is to place a prior distribution (often a Gaussian) over each parameter. The resulting model corresponds to a Gaussian process for infinite parameters (Neal, 1995), and a Bayesian NN (MacKay, 1992) for a finite number of parameters. Inference in BNNs is difficult however (Gal, 2016), so focus has thus shifted to techniques that approximate the posterior, *approximate BNNs*. Methods based on variational inference (VI) typically rely on a fully factorized approximate distribution (Kingma & Welling, 2014; Hinton & Van Camp, 1993), but often do not scale. To alleviate these difficulties, (Graves, 2011) proposed a model using sampling methods to estimate a factorized posterior. Probabilistic backpropagation (PBP), estimates a factorized posterior via expectation propagation (Hernández-Lobato & Adams, 2015).

Using several strategies to address scaling issues, Deep

Gaussian Processes show superior performance in terms of RMSE and uncertainty quality compared to state-of-the-art approximate BNNs (Bui et al., 2016)¹. Another recent approach to Bayesian learning, Bayesian hypernetworks, use a NN to learn a distribution of parameters over another network (Krueger et al., 2017). Multiplicative Normalizing Flows for variational Bayesian networks (MNF) (Louizos & Welling, 2017) is a recent model that formulates a posterior dependent on auxiliary variables. MNF achieves a highly flexible posterior by the application of normalizing flows to the auxiliary variables.

Although these recent techniques address some of the difficulties with approximate BNNs, they all require modifications to the architecture or the way networks are trained, as well as specialized knowledge from practitioners. Recently, (Gal & Ghahramani, 2015) showed that a network trained with dropout implicitly performs the VI objective. Therefore *any* network trained with dropout can be treated as an approximate Bayesian model by making multiple predictions through the network while sampling different dropout masks for each prediction. The mean and variance of the predictions are used in the estimation of the mean and variance of the predictive distribution².

3. Method

In the following, we introduce Bayesian models and a variational approximation using Kullback-Leibler (KL) divergence following (Gal, 2016). We continue by showing that a batch normalized deep network can be seen as an approximate Bayesian model. Employing theoretical insights and empirical analysis, we study the induced prior on the parameters when using batch normalization. Finally, we describe the procedure for estimating the uncertainty of a batch normalized network’s output.³

3.1. Bayesian Modeling

We assume a finite training set $\mathbf{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1:N}$ where each $(\mathbf{x}_i, \mathbf{y}_i)$ is a sample-label pair. Using \mathbf{D} , we are interested in learning an inference function $f_\omega(\mathbf{x}, \mathbf{y})$ with parameters ω . In deterministic models, the estimated label $\hat{\mathbf{y}}$ is obtained as follows:

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} f_\omega(\mathbf{x}, \mathbf{y})$$

In probabilistic models we let $f_\omega(\mathbf{x}, \mathbf{y}) = p(\mathbf{y}|\mathbf{x}, \omega)$. In Bayesian modeling, in contrast to finding a point estimate

¹By uncertainty quality, we refer to predictive probability distributions as measured by PLL and CRPS.

²This technique is referred to as “MC Dropout” in the original work, though we refer to it here as MCDO.

³While the method applies to FC or Conv layers, the induced prior from weight decay (Section 3.3) is studied for FC layers.

of the model parameters, the idea is to estimate an (approximate) posterior distribution of the model parameters $p(\boldsymbol{\omega}|\mathbf{D})$ to be used for probabilistic prediction:

$$p(\mathbf{y}|\mathbf{x}, \mathbf{D}) = \int f_{\boldsymbol{\omega}}(\mathbf{x}, \mathbf{y})p(\boldsymbol{\omega}|\mathbf{D})d\boldsymbol{\omega}$$

The predicted label, $\hat{\mathbf{y}}$, can then be accordingly obtained by sampling $p(\mathbf{y}|\mathbf{x}, \mathbf{D})$ or taking its maxima.

Variational Approximation In approximate Bayesian modeling, a common approach is to learn a parameterized approximating distribution $q_{\boldsymbol{\theta}}(\boldsymbol{\omega})$ that minimizes $\text{KL}(q_{\boldsymbol{\theta}}(\boldsymbol{\omega})||p(\boldsymbol{\omega}|\mathbf{D}))$; the Kullback-Leibler divergence of the true posterior w.r.t. its approximation. Minimizing this KL divergence is equivalent to the following minimization while being free of the data term $p(\mathbf{D})$ ⁴:

$$\begin{aligned} \mathcal{L}_{\text{VA}}(\boldsymbol{\theta}) := & - \sum_{i=1}^N \int q_{\boldsymbol{\theta}}(\boldsymbol{\omega}) \ln f_{\boldsymbol{\omega}}(\mathbf{x}_i, \mathbf{y}_i) d\boldsymbol{\omega} \\ & + \text{KL}(q_{\boldsymbol{\theta}}(\boldsymbol{\omega})||p(\boldsymbol{\omega})) \end{aligned}$$

During optimization, we want to take the derivative of the expected likelihood w.r.t. the learnable parameters $\boldsymbol{\theta}$. We use the same MC estimate as in (Gal, 2016) (explained in Appendix Section 1.1), such that one realized $\hat{\boldsymbol{\omega}}_i$ is taken for each sample i ⁵. Optimizing over mini-batches of size M , the approximated objective becomes:

$$\hat{\mathcal{L}}_{\text{VA}}(\boldsymbol{\theta}) := - \frac{N}{M} \sum_{i=1}^M \ln f_{\hat{\boldsymbol{\omega}}_i}(\mathbf{x}_i, \mathbf{y}_i) + \text{KL}(q_{\boldsymbol{\theta}}(\boldsymbol{\omega})||p(\boldsymbol{\omega})) \quad (1)$$

The first term is the data likelihood and the second term is the divergence of the prior w.r.t. the approximated posterior.

3.2. Batch Normalized Deep Nets as Bayesian Modeling

We now describe the optimization procedure of a deep network with batch normalization and draw the resemblance to the approximate Bayesian modeling in Eq (1).

The inference function of a feed-forward deep network with L layers can be described as:

$$f_{\boldsymbol{\omega}}(\mathbf{x}) = \mathbf{W}^L a(\mathbf{W}^{L-1} \dots a(\mathbf{W}^2 a(\mathbf{W}^1 \mathbf{x}))$$

⁴Achieved by constructing the Evidence Lower Bound, called ELBO, and assuming i.i.d. observation noise; details can be found in Appendix Section 1.1.

⁵While a MC integration using a single sample is a weak approximation, in an iterative optimization for $\boldsymbol{\theta}$ several samples will be taken over time.

where $a(\cdot)$ is an element-wise nonlinearity function and \mathbf{W}^l is the weight vector at layer l . Furthermore, we denote the input to layer l as \mathbf{x}^l with $\mathbf{x}^1 = \mathbf{x}$ and we then set $\mathbf{h}^l = \mathbf{W}^l \mathbf{x}^l$. Parenthesized super-index for matrices (e.g. $\mathbf{W}^{(j)}$) and vectors (e.g. $x^{(j)}$) indicates j th row and element respectively. Super-index u refers to a specific unit at layer l , (e.g. $\mathbf{W}^u = \mathbf{W}^{l,(j)}$, $h^u = h^{l,(j)}$).⁶

Batch Normalization Each layer of a deep network is constructed by several linear units whose parameters are the rows of the weight matrix \mathbf{W} . Batch normalization is a unit-wise operation proposed in (Ioffe & Szegedy, 2015) to standardize the distribution of each unit's input. For FC layers, it converts a unit's input h^u in the following way:

$$\hat{h}^u = \frac{h^u - \mathbb{E}[h^u]}{\sqrt{\text{Var}[h^u]}}$$

where the expectations are computed over the training set during evaluation, and mini-batch during training (in deep networks, the weight matrices are often optimized using back-propagated errors calculated on mini-batches of data)⁷. Therefore, during training, the estimated mean and variance on the mini-batch \mathbf{B} is used, which we denote by $\boldsymbol{\mu}_{\mathbf{B}}$ and $\boldsymbol{\sigma}_{\mathbf{B}}$ respectively. This makes the inference at training time for a sample \mathbf{x} a stochastic process, varying based on other samples in the mini-batch.

Loss Function and Optimization Training deep networks with mini-batch optimization involves a (regularized) risk minimization with the following form:

$$\mathcal{L}_{\text{RR}}(\boldsymbol{\omega}) := \frac{1}{M} \sum_{i=1}^M l(\hat{\mathbf{y}}_i, \mathbf{y}_i) + \Omega(\boldsymbol{\omega})$$

where the first term is the empirical loss on the training data and the second term is a regularization penalty acting as a prior on model parameters $\boldsymbol{\omega}$. If the loss l is cross-entropy for classification or sum-of-squares for regression problems (assuming i.i.d. Gaussian noise on labels), the first term is equivalent to minimizing the negative log-likelihood:

$$\mathcal{L}_{\text{RR}}(\boldsymbol{\omega}) := - \frac{1}{M\tau} \sum_{i=1}^M \ln f_{\boldsymbol{\omega}}(\mathbf{x}_i, \mathbf{y}_i) + \Omega(\boldsymbol{\omega})$$

⁶For a (softmax) classification network, $f_{\boldsymbol{\omega}}(\mathbf{x})$ is a vector with $f_{\boldsymbol{\omega}}(\mathbf{x}, \mathbf{y}) = f_{\boldsymbol{\omega}}(\mathbf{x})^{(\mathbf{y})}$, for regression networks with i.i.d. Gaussian noise we have $f_{\boldsymbol{\omega}}(\mathbf{x}, \mathbf{y}) = \mathcal{N}(f_{\boldsymbol{\omega}}(\mathbf{x}), \tau^{-1}\mathbf{I})$.

⁷It also learns an affine transformation for each unit with parameters $\boldsymbol{\gamma}$ and $\boldsymbol{\beta}$, omitted for brevity: $\hat{x}_{\text{affine}}^{(j)} = \boldsymbol{\gamma}^{(j)} \hat{x}^{(j)} + \boldsymbol{\beta}^{(j)}$.

with $\tau = 1$ for classification. In a network with batch normalization, the model parameters include $\{\mathbf{W}^{1:L}, \gamma^{1:L}, \beta^{1:L}, \mu_{\mathbf{B}}^{1:L}, \sigma_{\mathbf{B}}^{1:L}\}$. If we decouple the learnable parameters $\theta = \{\mathbf{W}^{1:L}, \gamma^{1:L}, \beta^{1:L}\}$ from the stochastic parameters $\omega = \{\mu_{\mathbf{B}}^{1:L}, \sigma_{\mathbf{B}}^{1:L}\}$, we get the following objective at each step of the mini-batch optimization:

$$\mathcal{L}_{\text{RR}}(\theta) := -\frac{1}{M\tau} \sum_{i=1}^M \ln f_{\{\theta, \hat{\omega}_i\}}(\mathbf{x}_i, \mathbf{y}_i) + \Omega(\theta) \quad (2)$$

where $\hat{\omega}_i$ is the means and variances for sample i 's mini-batch at a certain training step. Note that while $\hat{\omega}_i$ formally needs to be i.i.d. for each training example, a batch normalized network samples the stochastic parameters once per training step (mini-batch). For a large number of epochs, however, the distribution of sampled batch members for a given training example converges to the i.i.d. case.

In a batch normalized network, $q_{\theta}(\omega)$ corresponds to the joint distribution of the weights, induced by the randomness of the normalization parameters $\mu_{\mathbf{B}}^{1:L}, \sigma_{\mathbf{B}}^{1:L}$, as implied by the repeated sampling from \mathbf{D} during training. This is an approximation of the true posterior, where we have restricted the posterior to lie within the domain of our parametric network and source of randomness. With that, *we can estimate the uncertainty of predictions from a trained batch normalized network using the inherent stochasticity of BN* (Section 3.4).

3.3. Prior $p(\omega)$

Equivalence between the VA and BN training procedures requires $\frac{\partial}{\partial \theta}$ of Eq. (1) and Eq. (2) to be equivalent up to a scaling factor. This is the case if $\frac{\partial}{\partial \theta} \text{KL}(q_{\theta}(\omega) \| p(\omega)) = N\tau \frac{\partial}{\partial \theta} \Omega(\theta)$.

To reconcile this condition, one option is to let the prior $p(\omega)$ imply the regularization term $\Omega(\theta)$. Eq. (1) reveals that the contribution of $\text{KL}(q_{\theta}(\omega) \| p(\omega))$ to the optimization objective is inversely scaled with N . For BN, this corresponds to a model with a small $\Omega(\theta)$ when N is large. In the limit as $N \rightarrow \infty$, the optimization objectives of Eq. (1) and Eq. (2) become identical with no regularization.⁸

Another option is to let some $\Omega(\theta)$ imply $p(\omega)$. In Appendix Section 1.4 we explore this with L2-regularization, also called weight decay ($\Omega(\theta) = \lambda \sum_{l=1:L} \|\mathbf{W}^l\|^2$). We find that unlike in MCDO (Gal, 2016), some simplifying

⁸To prove the existence and find an expression of $\text{KL}(q_{\theta}(\omega) \| p(\omega))$, in Appendix Section 1.3 we find that BN approximately induces Gaussian distributions over BN units' means and standard deviations, centered around the population values given by \mathbf{D} . We assume a factorized distribution and Gaussian priors, and find the corresponding $\text{KL}(q_{\theta}(\omega) \| p(\omega))$ components in Appendix Section 1.4 Eq. (7). These could be used to construct a custom $\Omega(\theta)$ for any Gaussian choice of $p(\omega)$.

assumptions are necessary to reconcile the VA and BN objectives with weight decay: no scale and shift applied to BN layers, uncorrelated units in each layer, BN applied on all layers, and large N and M . Given these conditions:

$$\begin{aligned} p(\mu_{\mathbf{B}}^u) &= \mathcal{N}(\mu_{\mu,p}, \sigma_{\mu,p}) \\ p(\sigma_{\mathbf{B}}^u) &= \mathcal{N}(\mu_{\sigma,p}, \sigma_{\sigma,p}) \end{aligned}$$

where $\mu_{\mu,p} = 0$, $\sigma_{\mu,p} \rightarrow \infty$, $\mu_{\sigma,p} = 0$ and $\sigma_{\sigma,p} \rightarrow \frac{1}{2N\tau\lambda_l}$.

This corresponds to a wide and narrow distribution on BN units' means and std. devs respectively, where N accounts for the narrowness of the prior. Due to its popularity in deep learning, our experiments in Section 4 are performed with weight decay.

3.4. Predictive Uncertainty in Batch Normalized Deep Nets

In the absence of the true posterior, we rely on the approximate posterior to express an approximate predictive distribution:

$$p^*(\mathbf{y}|\mathbf{x}, \mathbf{D}) := \int f_{\omega}(\mathbf{x}, \mathbf{y}) q_{\theta}(\omega) d\omega$$

Following (Gal, 2016) we estimate the first (for regression and classification) and second (for regression) moments of the predictive distribution empirically (see Appendix Section 1.5 for details):

$$\begin{aligned} \mathbb{E}_{p^*}[\mathbf{y}] &\approx \frac{1}{T} \sum_{i=1}^T f_{\hat{\omega}_i}(\mathbf{x}) \\ \text{Cov}_{p^*}[\mathbf{y}] &\approx \tau^{-1} \mathbf{I} + \frac{1}{T} \sum_{i=1}^T f_{\hat{\omega}_i}(\mathbf{x})^{\top} f_{\hat{\omega}_i}(\mathbf{x}) \\ &\quad - \mathbb{E}_{p^*}[\mathbf{y}]^{\top} \mathbb{E}_{p^*}[\mathbf{y}] \end{aligned}$$

where each $\hat{\omega}_i$ corresponds to sampling the net's stochastic parameters $\omega = \{\mu_{\mathbf{B}}^{1:L}, \sigma_{\mathbf{B}}^{1:L}\}$ the same way as during training. Sampling $\hat{\omega}_i$ therefore involves sampling a batch \mathbf{B} from the *training set* and updating the parameters in the BN units, just as if we were taking a training step with \mathbf{B} . From a VA perspective, training the network amounted to minimizing $\text{KL}(q_{\theta}(\omega) \| p(\omega|\mathbf{D}))$ wrt θ . Sampling $\hat{\omega}_i$ from the training set, and keeping the size of \mathbf{B} consistent with the mini-batch size used during training, ensures that $q_{\theta}(\omega)$ during inference remains identical to the approximate posterior optimized during training.

The network is trained just as a regular BN network, but instead of replacing $\omega = \{\mu_{\mathbf{B}}^{1:L}, \sigma_{\mathbf{B}}^{1:L}\}$ with population values from \mathbf{D} for inference, we update these parameters stochastically, once for each forward pass.⁹ Pseudocode for estimating predictive mean and variance is given in Algorithm 1.

⁹As an alternative to using the training set \mathbf{D} to sample $\hat{\omega}_i$,

Algorithm 1 MCBN Algorithm

Input: sample x , number of inferences T , batchsize b
Output: mean prediction \hat{y} , predictive uncertainty σ^2

- 1: $\mathbf{y} = \{\}$
 - 2: **loop** for T iterations
 - 3: $B \sim D$ // mini batch
 - 4: $\hat{\omega} = \{\boldsymbol{\mu}_B, \boldsymbol{\sigma}_B\}$ // mini batch mean and variance
 - 5: $\mathbf{y} = \mathbf{y} \cup f_{\hat{\omega}}(x)$
 - 6: **end loop**
 - 7: $\hat{y} = \mathbb{E}[\mathbf{y}]$
 - 8: $\sigma^2 = \text{Cov}[\mathbf{y}] + \tau^{-1}\mathbf{I}$ // for regression
-

 $(\mathbf{y}_i, \mathbf{x}_i)$ as:

$$\text{PLL}(f_{\omega}(\mathbf{x}), (\mathbf{y}_i, \mathbf{x}_i)) = \log p(\mathbf{y}_i | f_{\omega}(\mathbf{x}_i))$$

where $p(\mathbf{y}_i | f_{\omega}(\mathbf{x}_i))$ is the model’s predicted PDF evaluated at \mathbf{y}_i , given the input x_i . A more detailed description is given in the Appendix Section 1.5. The metric is unbounded and maximized by a perfect prediction (mode at \mathbf{y}_i) with no variance. As the predictive mode moves away from \mathbf{y}_i , increasing the variance tends to increase PLL (by maximizing probability mass at \mathbf{y}_i). While PLL is an elegant measure, it has been criticized for allowing outliers to have an overly negative effect on the score (Selten, 1998).

4. Experiments and Results

We assess the uncertainty quality of MCBN quantitatively and qualitatively. Our quantitative analysis relies on CIFAR10 for image classification and eight standard regression datasets, listed in Appendix Table 1. Publicly available from the UCI Machine Learning Repository (University of California, 2017) and Delve (Ghahramani, 1996), these datasets have been used to benchmark comparative models in recent related literature (see (Hernández-Lobato & Adams, 2015), (Gal & Ghahramani, 2015), (Bui et al., 2016) and (Li & Gal, 2017)). We report results using standard metrics, and also propose useful upper and lower bounds to normalize these metrics for an easier interpretation in Section 4.2.

Our qualitative results include the toy dataset in Figure 1 in the style of (Karpathy, 2015), *a new visualization of uncertainty quality* that plots test errors sorted by predicted variance (Figure 2 and Appendix), and image segmentation results (Figure 2 and Appendix).

4.1. Metrics

We evaluate uncertainty quality based on two standard metrics, described below: Predictive Log Likelihood (PLL) and Continuous Ranked Probability Score (CRPS). To improve the interpretability of the metrics, we propose to normalize them by upper and lower bounds.

Predictive Log Likelihood (PLL) Predictive Log Likelihood is widely accepted as the main uncertainty quality metric for regression (Hernández-Lobato & Adams, 2015; Gal & Ghahramani, 2015; Bui et al., 2016; Li & Gal, 2017). A key property of PLL is that it makes no assumptions about the form of the distribution. The measure is defined for a probabilistic model $f_{\omega}(\mathbf{x})$ and a single observation

we could sample from the implied $q_{\theta}(\omega)$ as modeled in the Appendix. This would alleviate having to store \mathbf{D} for use during prediction. In our experiments we used \mathbf{D} to sample $\hat{\omega}_i$ however, and leave the evaluation of the modeled $q_{\theta}(\omega)$ for future research.

Continuous Ranked Probability Score (CRPS) Continuous Ranked Probability Score is a measure that takes the full predicted PDF into account with less sensitivity to outliers. A prediction with low variance that is slightly offset from the true observation will receive a higher score form CRPS than PLL. In order for CRPS to be analytically tractable, we need to assume a Gaussian unimodal predictive distribution. CRPS is defined as

$$\text{CRPS}(f_{\omega}(x_i), (y_i, x_i)) = \int_{-\infty}^{\infty} (F(y) - \mathbf{1}(y \geq y_i))^2 dy$$

where $F(y)$ is the predictive CDF, and $\mathbf{1}(y \geq y_i) = 1$ if $y \geq y_i$ and 0 otherwise (for univariate distributions) (Gneiting & Raftery, 2007). CRPS is interpreted as the sum of the squared area between the CDF and 0 where $y < y_i$ and between the CDF and 1 where $y \geq y_i$. A perfect prediction with no variance yields a CRPS of 0; for all other cases the value is larger. CRPS has no upper bound.

4.2. Benchmark models and normalized metrics

It is difficult to interpret the quality of uncertainty from raw PLL and CRPS values. We propose to normalize the metrics between useful lower and upper bounds. The normalized measures estimate the performance of an uncertainty model between the trivial solution (constant uncertainty) and *optimal uncertainty* for each prediction. For the lower bound, we define a baseline that predicts constant variance regardless of input. The variance is set to a fixed value that optimizes CRPS on validation data. We call this model Constant Uncertainty BN (CUBN). It reflects our best guess of constant variance on test data – thus, any improvement in uncertainty quality over CUBN indicates a sensible estimate of uncertainty. We similarly define a baseline for dropout, Constant Uncertainty Dropout (CUDO). The modeling of variance (uncertainty) by MCBN and CUBN are visualized in Figure 1.

An upper bound on uncertainty performance can also be defined for a probabilistic model f with respect to CRPS or PLL. For each observation (y_i, x_i) , a value

for the predictive variance T_i can be chosen that maximizes PLL or minimizes CRPS¹⁰. Using CUBN as a lower bound and the optimized CRPS score as the upper bound, uncertainty estimates can be normalized between these bounds (1 indicating optimal performance, and 0 indicating same performance as fixed uncertainty). We call this normalized measure $\overline{\text{CRPS}} = \frac{\text{CRPS}(f, (y_i, x_i)) - \text{CRPS}(f_{\text{CU}}, (y_i, x_i))}{\min_T \text{CRPS}(f, (y_i, x_i)) - \text{CRPS}(f_{\text{CU}}, (y_i, x_i))} \times 100$, and the PLL analogue $\overline{\text{PLL}} = \frac{\text{PLL}(f, (y_i, x_i)) - \text{PLL}(f_{\text{CU}}, (y_i, x_i))}{\max_T \text{PLL}(f, (y_i, x_i)) - \text{PLL}(f_{\text{CU}}, (y_i, x_i))} \times 100$.

4.3. Test setup

Our evaluation compares MCBN to MCDO (Gal & Ghahramani, 2015) and MNF (Louizos & Welling, 2017) using the datasets and metrics described above. Our setup is similar to (Hernández-Lobato & Adams, 2015), which was also followed by (Gal & Ghahramani, 2015). However, our comparison implements a different hyperparameter selection, allows for a larger range of dropout rates, and uses larger networks with two hidden layers.

For the regression task, all models share a similar architecture: two hidden layers with 50 units each, and ReLU activations, with the exception of Protein Tertiary Structure dataset (100 units per hidden layer). Inputs and outputs were normalized during training. Results were averaged over five random splits of 20% test and 80% training and cross-validation (CV) data. For each split, 5-fold CV by grid search with a RMSE minimization objective was used to find training hyperparameters and optimal n.o. epochs, out of a maximum of 2000. For BN-based models, the hyperparameter grid consisted of a weight decay factor ranging from 0.1 to 1^{-15} by a log 10 scale, and a batch size range from 32 to 1024 by a log 2 scale. For DO-based models, the hyperparameter grid consisted of the same weight decay range, and dropout probabilities in $\{0.2, 0.1, 0.05, 0.01, 0.005, 0.001\}$. DO-based models used a batch size of 32 in all evaluations. For MNF¹¹, the n.o. epochs was optimized, the batch size was set to 100, and early stopping test performed each epoch (compared to every 20th for MCBN, MCDO).

For MCBN and MCDO, the model with optimal training hyperparameters was used to optimize τ numerically. This optimization was made in terms of average CV CRPS for MCBN, CUBN, MCDO, and CUDO respectively.

Estimates for the predictive distribution were obtained by taking $T = 500$ stochastic forward passes through the network. For each split, test set evaluation was done 5 times with different seeds. Implementation was done in TensorFlow with the Adam optimizer and a learning rate of 0.001.

¹⁰ T_i can be found analytically for PLL, but must be found numerically for CRPS.

¹¹Where we used an adapted version of the authors' code.

For the image classification test we use CIFAR10 (Krizhevsky & Hinton, 2009) which includes 10 object classes with 5,000 and 1,000 images in the training and test sets, respectively. Images are 32x32 RGB format. We trained a ResNet32 architecture with a batch size of 32, learning rate of 0.1, weight decay of 0.0002, leaky ReLU slope of 0.1, and 5 residual units. SGD with momentum was used as the optimizer.

Code for reproducing our experiments is available at <https://github.com/icml-mcbsn/mcbsn>.

4.4. Test results

The regression experiment comparing uncertainty quality is summarized in Table 1. We report $\overline{\text{CRPS}}$ and $\overline{\text{PLL}}$, expressed as a percentage, which reflects how close the model is to the upper bound, and check to see if the model significantly exceeds the lower bound using a one sample t-test (significance level is indicated by *'s). Further details are provided in Appendix Section 1.7.

In Figure 2 (left), we present a novel visualization of uncertainty quality for regression problems. Data are sorted by estimated uncertainty in the x -axis. Grey dots show the errors in model predictions, and the shaded areas show the model uncertainty. A running mean of the errors appears as a gray line. If uncertainty estimation is working well, a correlation should exist between the mean error (gray line) and uncertainty (shaded area). This indicates that the uncertainty estimation recognizes samples with larger (or smaller) potential for predictive errors.

We applied MCBN on the image classification task of CIFAR10. The baseline in this case is the softmax distribution using the moving average for BN units. Log likelihood (PLL) is the metric used to compare with the baseline. The baseline achieves a PLL of **-0.32** on the test set, while MCBN obtains a PLL of **-0.28**. Table 2 shows the performance of MCBN when using different number of stochastic forward passes (the MCBN batchsize is fixed to the training batch size at 32). PLL improves as the number of the stochastic passes increases, until it is significantly better than the softmax baseline.

To demonstrate how model uncertainty can be obtained from an existing network with minimal effort, we applied MCBN to an image segmentation task using Bayesian SegNet with the main CamVid and PASCAL-VOC models in (Kendall et al., 2015). We simply ran multiple forward passes with different mini-batches randomly taken from the train set. The models obtained from the online model zoo have BN blocks after each layer. We recalculate mean and variance for the first 2 blocks only and use the training statistics for the rest of the blocks. Mini-batches of size 10 and 36 were used for CamVid and VOC respectively

Table 1. **Uncertainty quality measured on eight regression datasets.** MCBN, MCDO and MNF are compared over 5 random 80-20 splits of the data with 5 different random seeds each split. We report $\overline{\text{CRPS}}$ and $\overline{\text{PLL}}$, uncertainty metrics CRPS and PLL normalized to a lower bound of constant variance and upper bound that maximizes the metric expressed as a percentage (described in Section 4.2). Higher numbers mean the model is closer to the upper bound. We check if the reported values for CRPS and $\overline{\text{PLL}}$ significantly exceed the lower bound using a one sample t-test (significance level indicated by *’s). See text for further details.

Dataset	$\overline{\text{CRPS}}$			$\overline{\text{PLL}}$		
	MCBN	MCDO	MNF	MCBN	MCDO	MNF
Boston	8.50 ****	3.06 ****	5.88 ****	10.49 ****	5.51 ****	1.76 ns
Concrete	3.91 ****	0.93 *	3.13 ***	-36.36 **	10.92 ****	-2.16 ns
Energy	5.75 ****	1.37 ns	1.10 ns	10.89 ****	-14.28 *	-33.88 ns
Kin8nm	2.85 ****	1.82 ****	0.53 ns	1.68 ***	-0.26 ns	0.42 ns
Power	0.24 ***	-0.44 ****	-0.89 ****	0.33 **	3.52 ****	-0.87 ****
Protein	2.66 ****	0.99 ****	0.57 ****	2.56 ****	6.23 ****	0.52 ****
Wine (Red)	0.26 **	2.00 ****	0.94 ****	0.19 *	2.91 ****	0.83 ****
Yacht	-56.39 ***	21.42 ****	24.92 ****	45.58 ****	-41.54 ns	46.19 ****

due to memory limits. The results in Figure 2 (*right*) were obtained from 20 stochastic forward passes, showing high uncertainty near object boundaries. The VOC results are more appealing because of larger mini-batches.

We provide additional experimental results in the Appendix. Appendix Tables 2 and 3 show the mean $\overline{\text{CRPS}}$ and $\overline{\text{PLL}}$ values from the regression experiment. Table 4 provides the raw CRPS and PLL scores. In Table 5 we provide RMSE results of the MCBN and MCDO networks in comparison with non-stochastic BN and DO networks. These results indicate that the procedure of multiple forward passes in MCBN and MCDO show slight improvements in the predictive accuracy compared to their non-Bayesian counterparts. In Tables 6 and 7, we investigate the effect of varying batch size while keeping other hyperparameters fixed. We see that performance deteriorates with small batch sizes (≤ 16), a known issue of BN (Ioffe, 2017). Similarly, results varying the number of stochastic forward passes T is reported in Tables 8 and 9. While performance benefits from large T , in some cases $T = 50$ (i.e. 1/10 of T in the main evaluation) performs well. Uncertainty-error plots for all the datasets are provided in the Appendix.

5. Discussion

The results presented in Tables 1-2 and Appendix Tables 2-9 indicate that MCBN generates meaningful uncertainty

Table 2. **Uncertainty quality for image classification varying number of stochastic forward passes.** Uncertainty quality for image classification measured by PLL. ResNet32 is trained on CIFAR10 with batch size 32. PLL improves as the sampling increases until it is significantly better than the softmax baseline (-0.32).

	Number of stochastic forward passes							
	1	2	4	8	16	32	64	128
PLL	-.36	-.32	-.30	-.29	-.29	-.28	-.28	-.28

estimates that correlate with actual errors in the model’s prediction. In Table 1, we show statistically significant improvements over CUBN in the majority of the datasets, both in terms of $\overline{\text{CRPS}}$ and $\overline{\text{PLL}}$. The visualizations in Figure 2 and in the Appendix Figures 2-3 show correlations between the estimated model uncertainty and errors of the network’s predictions. We perform the same experiments using MCDO and MNF, and find that MCBN generally performs on par with both methods. Looking closer, MCBN outperforms MCDO and MNF in more cases than not, measured by $\overline{\text{CRPS}}$. However, care must be used. The learned parameters are different, leading to different predictive means and confounding direct comparison.

The results on the Yacht Hydrodynamics dataset seem contradictory. The $\overline{\text{CRPS}}$ score for MCBN are extremely negative, while the $\overline{\text{PLL}}$ score is extremely positive. The opposite trend is observed for MCDO. To add to the puzzle, the visualization in Figure 2 depicts an extremely promising uncertainty estimation that models the predictive errors with high fidelity. We hypothesize that this strange behavior is due to the small size of the data set, which only contains 60 test samples, or due to the Gaussian assumption of CRPS. There is also a large variability in the model’s accuracy on this dataset, which further confounds the measurements for such limited data.

One might criticize the overall quality of uncertainty estimates observed in all the models we tested, due to the magnitude of $\overline{\text{CRPS}}$ and $\overline{\text{PLL}}$ in Table 1. The scores rarely exceed 10% improvement over the lower bound. However, we caution that these measures should be taken in context. The upper bound is very difficult to achieve in practice – it is optimized for *each test sample individually* – and the lower bound is a quite reasonable estimate.

The study of MCBN sensitivity to batch size revealed that a certain batch size is required for the best performance, dependent on the data. When doing inference on a GPU, large

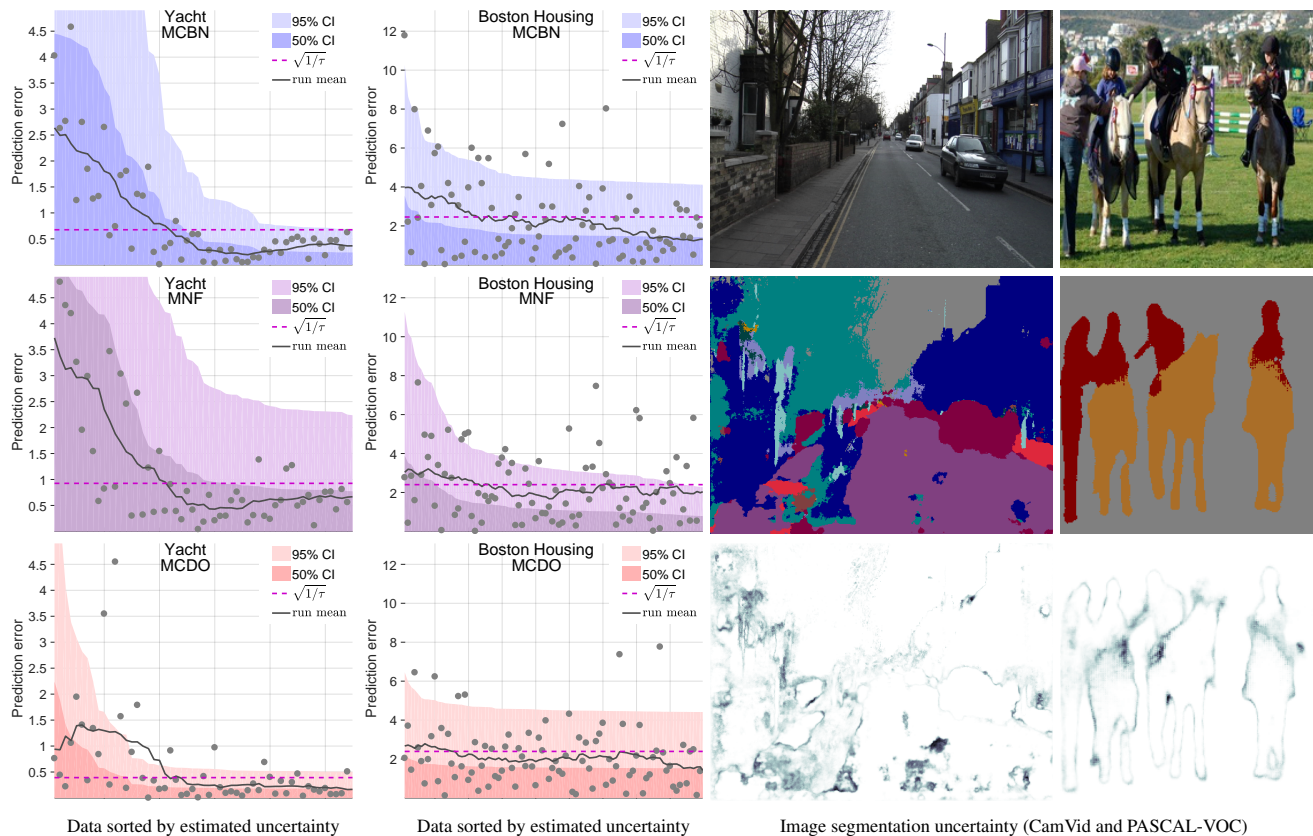


Figure 2. Uncertainty-error plots (left) and segmentation and uncertainty results applying MCBN to Bayesian SegNet (right). (left) Errors in predictions (gray dots) sorted by estimated uncertainty on select datasets. The shaded areas show model uncertainty for MCBN (blue), MNF (violet) and MCDO (red). The light area indicates 95% CI, dark area 50% CI. Gray dots show absolute prediction errors on the test set, and the gray line depicts a running mean of the errors. The dashed line indicates the optimized constant uncertainty. A correlation between estimated uncertainty (shaded area) and mean error (gray) indicates the uncertainty estimates are meaningful for estimating errors. (right) Applying MCBN to Bayesian SegNet (Kendall et al., 2015) on scenes from CamVid (3rd column) and PASCAL-VOC (4th column). Top: original image. Middle: the Bayesian estimated segmentation. Bottom: estimated uncertainty using MCBN for all classes. The uncertainty maps for both datasets are reasonable, but qualitatively better for PASCAL-VOC due to the larger mini-batch size (36) compared to CamVid (10). Smaller batch sizes were used for CamVid due to memory limitations (CamVid images are 360x480 while VOC are 224x224). See Appendix for complete results.

batch sizes may cause memory issues for cases where the input is large and the network has a large number of parameters, as is common for state-of-the-art image classification networks. However, there are various workarounds to this problem. One can store BN statistics, instead of batches, to reduce memory issues. Furthermore, we can use the Gaussian estimate of the BN statistics as discussed previously, which makes memory and computation extremely efficient.

6. Conclusion

In this work, we have shown that training a deep network using batch normalization is equivalent to approximate inference in Bayesian models. We show evidence that the uncertainty estimates from MCBN correlate with actual errors in the model’s prediction, and are useful for practical

tasks such as regression, image classification, and image segmentation. Our experiments show that MCBN yields a significant improvement over the optimized constant uncertainty baseline, on par with MCDO and MNF. Our evaluation also suggests new normalized metrics based on useful upper and lower bounds, and a new visualization which provides an intuitive explanation of uncertainty quality.

Finally, it should be noted that over the past few years, batch normalization has become an integral part of most – if not all – cutting edge deep networks. We have shown that it is possible to obtain meaningful uncertainty estimates from existing models without modifying the network or the training procedure. With a few lines of code, robust uncertainty estimates can be obtained by computing the variance of multiple stochastic forward passes through an existing network.

References

- Bui, T. D., Hernández-Lobato, D., Li, Y., Hernández-Lobato, J. M., and Turner, R. E. Deep Gaussian Processes for Regression using Approximate Expectation Propagation. In *ICML*, 2016.
- Chen, X., Kundu, K., Zhang, Z., Ma, H., Fidler, S., and Urtasun, R. Monocular 3d object detection for autonomous driving. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2147–2156, 2016.
- Djuric, U., Zadeh, G., Aldape, K., and Diamandis, P. Precision histology: how deep learning is poised to revitalize histomorphology for personalized cancer care. *npj Precision Oncology*, 1(1):22, 2017.
- Esteva, A., Kuprel, B., Novoa, R. A., Ko, J., Swetter, S. M., Blau, H. M., and Thrun, S. Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, Feb 2017.
- Gal, Y. *Uncertainty in Deep Learning*. PhD thesis, University of Cambridge, 2016.
- Gal, Y. and Ghahramani, Z. Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning. *ICML*, 48:1–10, 2015.
- Ghahramani, Z. *Delve Datasets*. University of Toronto, 1996. URL <http://www.cs.toronto.edu/~delve/data/kin/desc.html>.
- Ghahramani, Z. Probabilistic machine learning and artificial intelligence. *Nature*, 521(7553):452–459, May 2015.
- Gneiting, T. and Raftery, A. E. Strictly Proper Scoring Rules, Prediction, and Estimation. *Journal of the American Statistical Association*, 102(477):359–378, 2007.
- Goodfellow, I. J., Shlens, J., and Szegedy, C. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- Graves, A. Practical Variational Inference for Neural Networks. *NIPS*, 2011.
- Hernández-Lobato, J. M. and Adams, R. Probabilistic backpropagation for scalable learning of bayesian neural networks. In *International Conference on Machine Learning*, pp. 1861–1869, 2015.
- Hinton, G. E. and Van Camp, D. Keeping the neural networks simple by minimizing the description length of the weights. In *Proceedings of the sixth annual conference on Computational learning theory*, pp. 5–13. ACM, 1993.
- Ioffe, S. Batch renormalization: Towards reducing mini-batch dependence in batch-normalized models. *CoRR*, abs/1702.03275, 2017. URL <http://arxiv.org/abs/1702.03275>.
- Ioffe, S. and Szegedy, C. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. *Arxiv*, 2015. URL <http://arxiv.org/abs/1502.03167>.
- Karpathy, A. Convnetjs demo: toy 1d regression, 2015. URL <http://cs.stanford.edu/people/karpathy/convnetjs/demo/regression.html>.
- Kendall, A., Badrinarayanan, V., and Cipolla, R. Bayesian SegNet: Model Uncertainty in Deep Convolutional Encoder-Decoder Architectures for Scene Understanding. *CoRR*, abs/1511.0, 2015. URL <http://arxiv.org/abs/1511.02680>.
- Kingma, D. P. and Welling, M. Auto-Encoding Variational Bayes. In *ICLR*, 2014.
- Krizhevsky, A. and Hinton, G. Learning multiple layers of features from tiny images. 2009.
- Krueger, D., Huang, C.-W., Islam, R., Turner, R., Lacoste, A., and Courville, A. Bayesian hypernetworks. *arXiv preprint arXiv:1710.04759*, 2017.
- Lehmann, E. L. *Elements of Large-Sample Theory*. Springer Verlag, New York, 1999. ISBN 0387985956.
- Li, Y. and Gal, Y. Dropout Inference in Bayesian Neural Networks with Alpha-divergences. *arXiv*, 2017.
- Louizos, C. and Welling, M. Multiplicative normalizing flows for variational Bayesian neural networks. In Precup, D. and Teh, Y. W. (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 2218–2227, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR. URL <http://proceedings.mlr.press/v70/louizos17a.html>.
- MacKay, D. J. A practical bayesian framework for back-propagation networks. *Neural computation*, 4(3):448–472, 1992.
- Neal, R. M. *Bayesian Learning for Neural Networks*. PhD thesis, University of Toronto, 1995.
- Neal, R. M. *Bayesian learning for neural networks*, volume 118. Springer Science & Business Media, 2012.
- Selten, R. Axiomatic characterization of the quadratic scoring rule. *Experimental Economics*, 1(1):43–62, 1998.

Shen, L. End-to-end training for whole image breast cancer diagnosis using an all convolutional design. *arXiv preprint arXiv:1708.09427*, 2017.

University of California, I. UC Irvine Machine Learning Repository, 2017. URL <https://archive.ics.uci.edu/ml/index.html>.

Wang, S. I. and Manning, C. D. Fast dropout training. *Proceedings of the 30th International Conference on Machine Learning*, 28:118–126, 2013. URL <http://machinelearning.wustl.edu/mlpapers/papers/wang13a>.

1. Appendix

1.1. Variational Approximation

Assume we were to come up with a family of distributions parameterized by θ in order to approximate the posterior, $q_\theta(\omega)$. Our goal is to set θ such that $q_\theta(\omega)$ is as similar to the true posterior $p(\omega|\mathbf{D})$ as possible.

For clarity, $q_\theta(\omega)$ is a distribution over stochastic parameters ω that is determined by a set of learnable parameters θ and some source of randomness. The approximation is therefore limited by our choice of parametric function $q_\theta(\omega)$ as well as the randomness.¹² Given ω and an input \mathbf{x} , an output distribution $p(\mathbf{y}|\mathbf{x}, \omega) = p(\mathbf{y}|f_\omega(\mathbf{x})) = f_\omega(\mathbf{x}, \mathbf{y})$ is induced by observation noise (the conditionality of which is omitted for brevity).

One strategy for optimizing θ is to minimize $\text{KL}(q_\theta(\omega)||p(\omega|\mathbf{D}))$, the KL divergence of $p(\omega|\mathbf{D})$ w.r.t. $q_\theta(\omega)$. Minimizing $\text{KL}(q_\theta(\omega)||p(\omega|\mathbf{D}))$ is equivalent to maximizing the ELBO:

$$\int_{\omega} q_\theta(\omega) \ln p(\mathbf{Y}|\mathbf{X}, \omega) d\omega - \text{KL}(q_\theta(\omega)||p(\omega))$$

Assuming i.i.d. observation noise, this is equivalent to minimizing:

$$\mathcal{L}_{\text{VA}}(\theta) := - \sum_{n=1}^N \int q_\theta(\omega) \ln p(\mathbf{y}_i|f_\omega(\mathbf{x}_i)) d\omega + \text{KL}(q_\theta(\omega)||p(\omega))$$

Instead of making the optimization on the full training set, we can use a subsampling (yielding an unbiased estimate of $\mathcal{L}_{\text{VA}}(\theta)$) for iterative optimization (as in mini-batch optimization):

$$\hat{\mathcal{L}}_{\text{VA}}(\theta) := - \frac{N}{M} \sum_{i \in B} \int_{\omega} q_\theta(\omega) \ln p(\mathbf{y}_i|f_\omega(\mathbf{x}_i)) d\omega + \text{KL}(q_\theta(\omega)||p(\omega))$$

During optimization, we want to take the derivative of the expected log likelihood w.r.t. the learnable parameters θ . (Gal, 2016) provides an intuitive interpretation of a MC estimate for NNs trained with a SRT (equivalent to the reparametrisation trick in (Kingma & Welling, 2014)), and this interpretation is followed here. We let an auxiliary variable ϵ represent the stochasticity during training such that $\omega = g(\theta, \epsilon)$. The function g and the distribution of ϵ are such that $p(g(\theta, \epsilon)) = q_\theta(\omega)$.¹³ Assuming $q_\theta(\omega)$ can be written $\int_{\epsilon} q_\theta(\omega|\epsilon) p(\epsilon) d\epsilon$ where $q_\theta(\omega|\epsilon) = \delta(\omega - g(\theta, \epsilon))$, this auxiliary variable yields the estimate (full proof in (Gal, 2016)):

$$\hat{\mathcal{L}}_{\text{VA}}(\theta) = - \frac{N}{M} \sum_{i \in B} \int_{\epsilon} p(\epsilon) \ln p(\mathbf{y}_i|f_{g(\theta, \epsilon)}(\mathbf{x}_i)) d\epsilon + \text{KL}(q_\theta(\omega)||p(\omega))$$

where taking a single sample MC estimate of the integral yields the optimization objective in Eq. 1.

¹²In an approx. Bayesian view of a NN, $q_\theta(\omega)$ would correspond to the distribution of weights in the network given by some SRT.

¹³In a NN trained with BN, it is easy to see that g exists if we let ϵ represent a mini-batch selection from the training data, since all BN units' means and variances are completely determined by ϵ and θ .

1.2. KL Divergence of factorized Gaussians

If $q_{\theta}(\omega)$ and $p(\omega)$ factorize over all stochastic parameters:

$$\begin{aligned}
 \text{KL}(q_{\theta}(\omega)||p(\omega)) &= - \int_{\omega} \prod_i [q_{\theta}(\omega_i)] \ln \frac{\prod_i p(\omega_i)}{\prod_i q_{\theta}(\omega_i)} d\omega \\
 &= - \int_{\omega} \prod_i [q_{\theta}(\omega_i)] \sum_i \left[\ln \frac{p(\omega_i)}{q_{\theta}(\omega_i)} \right] \prod_i d\omega_i \\
 &= \sum_j \left[- \int_{\omega} \prod_i [q_{\theta}(\omega_i)] \ln \frac{p(\omega_j)}{q_{\theta}(\omega_j)} \prod_i d\omega_i \right] \\
 &= \sum_j \left[- \int_{\omega_j} q_{\theta}(\omega_j) \ln \frac{p(\omega_j)}{q_{\theta}(\omega_j)} d\omega_j \int_{\omega_{i \neq j}} \prod_i q_{\theta}(\omega_i) d\omega_i \right] \\
 &= \sum_i - \int_{\omega_i} q_{\theta}(\omega_i) \ln \frac{p(\omega_i)}{q_{\theta}(\omega_i)} d\omega_i \\
 &= \sum_i \text{KL}(q_{\theta}(\omega_i)||p(\omega_i))
 \end{aligned} \tag{3}$$

such that $\text{KL}(q_{\theta}(\omega)||p(\omega))$ is the sum of the KL divergence terms for the individual stochastic parameters ω_i . If the factorized distributions are Gaussians, where $q_{\theta}(\omega_i) = \mathcal{N}(\mu_q, \sigma_q^2)$ and $p(\omega_i) = \mathcal{N}(\mu_p, \sigma_p^2)$ we get:

$$\begin{aligned}
 \text{KL}(q_{\theta}(\omega_i)||p(\omega_i)) &= \int_{\omega_i} q_{\theta}(\omega_i) \ln \frac{q_{\theta}(\omega_i)}{p(\omega_i)} d\omega_i \\
 &= -H(q_{\theta}(\omega_i)) - \int_{\omega_i} q_{\theta}(\omega_i) \ln p(\omega_i) d\omega_i \\
 &= -\frac{1}{2}(1 + \ln(2\pi\sigma_q^2)) \\
 &\quad - \int_{\omega_i} q_{\theta}(\omega_i) \ln \frac{1}{(2\pi\sigma_p^2)^{1/2}} \exp \left\{ -\frac{(\omega_i - \mu_p)^2}{2\sigma_p^2} \right\} d\omega_i \\
 &= -\frac{1}{2}(1 + \ln(2\pi\sigma_q^2)) \\
 &\quad + \frac{1}{2} \ln(2\pi\sigma_p^2) + \frac{\mathbb{E}_q[\omega_i^2] - 2\mu_p \mathbb{E}_q[\omega_i] + \mu_p^2}{2\sigma_p^2} \\
 &= \ln \frac{\sigma_p}{\sigma_q} + \frac{\sigma_q^2 + (\mu_q - \mu_p)^2}{2\sigma_p^2} - \frac{1}{2}
 \end{aligned} \tag{4}$$

for each KL divergence term. Here $H(q_{\theta}(\omega_i)) = \frac{1}{2}(1 + \ln(2\pi\sigma_q^2))$ is the differential entropy of $q_{\theta}(\omega_i)$.

1.3. Distribution of $\mu_{\mathbf{B}}^u, \sigma_{\mathbf{B}}^u$

Here we approximate the distribution of mean and standard deviation of a mini-batch, separately to two Gaussians – This has also been empirically verified, see Figure 3 for 2 sample plots and the appendix section 1.9 for more. For the mean we get:

$$\mu_{\mathbf{B}} = \frac{\sum_{m=1}^M \mathbf{W}^{(j)} \mathbf{x}_m}{M}$$

where \mathbf{x}_m are the examples in the sampled batch. We will assume these are sampled i.i.d.¹⁴. Samples of the random variable $\mathbf{W}^{(j)} \mathbf{x}_m$ are then i.i.d.. Then by central limit theorem (CLT) the following holds for sufficiently large M (often ≥ 30):

$$\mu_{\mathbf{B}} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{M}\right)$$

¹⁴Although in practice with deep learning, mini-batches are sampled without replacement, stochastic gradient descent samples with replacement in its standard form.

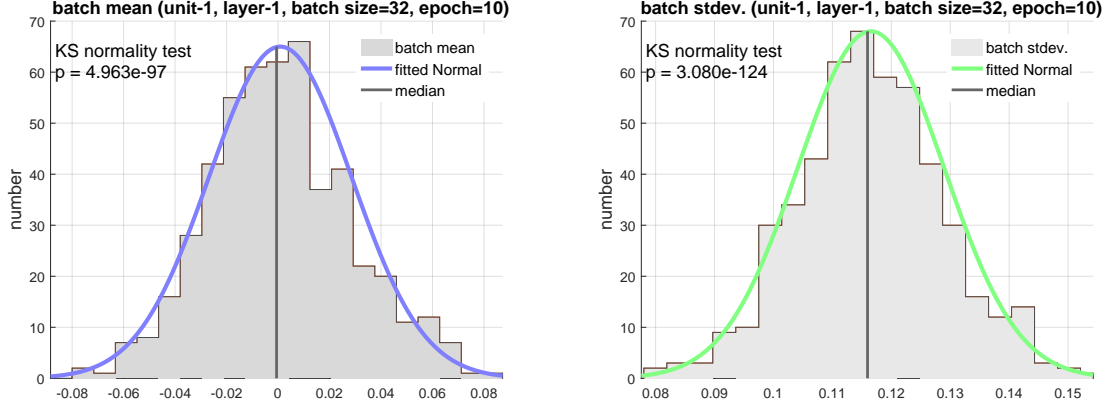


Figure 3. Batch statistics used to train the network are normal. A one-sample Kolmogorov-Smirnov test checks that μ_B and σ_B come from a standard normal distribution. More examples are available in Appendix 1.9.

For standard deviation:

$$\sigma_B = \sqrt{\frac{\sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu_B)^2}{M}}$$

Then

$$\sqrt{M}(\sigma_B - \sigma) = \sqrt{M} \left(\sqrt{\frac{\sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu_B)^2}{M}} - \sqrt{\sigma^2} \right)$$

We want to rewrite $\sqrt{\frac{\sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu_B)^2}{M}}$. We take a Taylor expansion of $f(x) = \sqrt{x}$ around $a = \sigma^2$. With $x = \frac{\sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu_B)^2}{M}$:

$$\sqrt{x} = \sqrt{\sigma^2} + \frac{1}{2\sqrt{\sigma^2}}(x - \sigma^2) + \mathcal{O}[(x - \sigma^2)^2]$$

so

$$\begin{aligned} \sqrt{M}(\sigma_B - \sigma) &= \sqrt{M} \left(\frac{1}{2\sqrt{\sigma^2}} \left(\frac{\sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu_B)^2}{M} - \sigma^2 \right) + \right. \\ &\quad \left. \mathcal{O} \left[\left(\frac{\sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu_B)^2}{M} - \sigma^2 \right)^2 \right] \right) \\ &= \frac{\sqrt{M}}{2\sigma} \left(\frac{1}{M} \sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu_B)^2 - \sigma^2 \right) + \\ &\quad \mathcal{O} \left[\sqrt{M} \left(\frac{\sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu_B)^2}{M} - \sigma^2 \right)^2 \right] \\ &= \frac{1}{2\sigma\sqrt{M}} \left(\sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu_B)^2 - M\sigma^2 \right) + \\ &\quad \mathcal{O} \left[\sqrt{M} \left(\frac{\sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu_B)^2}{M} - \sigma^2 \right)^2 \right] \end{aligned}$$

consider $\sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu_B)^2$. We know that $E[\mathbf{W}^{(j)} \mathbf{x}_m] = \mu$ and write

$$\begin{aligned}
 & \sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu_B)^2 \\
 &= \sum_{m=1}^M ((\mathbf{W}^{(j)} \mathbf{x}_m - \mu) - (\mu_B - \mu))^2 \\
 &= \sum_{m=1}^M ((\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^2 + (\mu_B - \mu)^2 - 2(\mathbf{W}^{(j)} \mathbf{x}_m - \mu)(\mu_B - \mu)) \\
 &= \sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^2 + M(\mu_B - \mu)^2 - 2(\mu_B - \mu) \sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu) \\
 &= \sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^2 - M(\mu_B - \mu)^2 \\
 &= \sum_{m=1}^M ((\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^2 - (\mu_B - \mu)^2)
 \end{aligned}$$

then

$$\begin{aligned}
 \sqrt{M}(\sigma_B - \sigma) &= \frac{1}{2\sigma\sqrt{M}} \left(\sum_{m=1}^M ((\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^2 - (\mu_B - \mu)^2) - M\sigma^2 \right) + \\
 & \quad \mathcal{O} \left[\sqrt{M} \left(\frac{\sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu_B)^2}{M} - \sigma^2 \right)^2 \right] \\
 &= \frac{1}{2\sigma\sqrt{M}} \left(\sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^2 - \sum_{m=1}^M (\mu_B - \mu)^2 - M\sigma^2 \right) + \\
 & \quad \mathcal{O} \left[\sqrt{M} \left(\frac{\sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu_B)^2}{M} - \sigma^2 \right)^2 \right] \\
 &= \frac{1}{2\sigma\sqrt{M}} \left(\sum_{m=1}^M ((\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^2 - \sigma^2) - \sum_{m=1}^M (\mu_B - \mu)^2 \right) + \\
 & \quad \mathcal{O} \left[\sqrt{M} \left(\frac{\sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu_B)^2}{M} - \sigma^2 \right)^2 \right] \\
 &= \frac{1}{2\sigma\sqrt{M}} \sum_{m=1}^M ((\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^2 - \sigma^2) \\
 & \quad - \frac{1}{2\sigma\sqrt{M}} \sum_{m=1}^M (\mu_B - \mu)^2 \\
 & \quad + \mathcal{O} \left[\sqrt{M} \left(\frac{\sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu_B)^2}{M} - \sigma^2 \right)^2 \right] \\
 &= \underbrace{\frac{1}{2\sigma\sqrt{M}} \sum_{m=1}^M ((\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^2 - \sigma^2)}_{\text{term A}} \\
 & \quad - \underbrace{\frac{\sqrt{M}}{2\sigma} (\mu_B - \mu)^2}_{\text{term B}} \\
 & \quad + \underbrace{\mathcal{O} \left[\sqrt{M} \left(\frac{\sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu_B)^2}{M} - \sigma^2 \right)^2 \right]}_{\text{term C}}
 \end{aligned}$$

We go through each term in turn

Term A

We have

$$\text{Term A} = \frac{1}{2\sigma\sqrt{M}} \sum_{m=1}^M ((\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^2 - \sigma^2)$$

where $\sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^2$ is the sum of M RVs $(\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^2$. Note that since $E[\mathbf{W}^{(j)} \mathbf{x}_m] = \mu$ it holds that $E[(\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^2] = \sigma^2$. Since $(\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^2$ is sampled approximately iid (by assumptions above), for large enough

M by CLT it holds approximately that

$$\sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^2 \sim \mathcal{N}(M\sigma^2, M\text{Var}((\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^2))$$

where

$$\begin{aligned} \text{Var}((\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^2) &= E[(\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^{2*2}] - E[(\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^2]^2 \\ &= E[(\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^4] - \sigma^4 \end{aligned}$$

Then

$$\sum_{m=1}^M ((\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^2 - \sigma^2) \sim \mathcal{N}(0, M * E[(\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^4] - M\sigma^4)$$

so

$$\text{Term A} \sim \mathcal{N}\left(0, \frac{E[(\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^4] - \sigma^4}{4\sigma^2}\right)$$

Term B

We have

$$\text{Term B} = \frac{\sqrt{M}}{2\sigma} (\mu_B - \mu)^2 = \frac{1}{2\sigma} \sqrt{M} (\mu_B - \mu) (\mu_B - \mu)$$

Consider $(\mu_B - \mu)$. As $\mu_B \xrightarrow{p} \mu$ when $M \rightarrow \infty$ we have $\mu_B - \mu \xrightarrow{p} 0$. We also have

$$\sqrt{M}(\mu_B - \mu) = \frac{\sum_{m=1}^M \mathbf{W}^{(j)} \mathbf{x}_m}{\sqrt{M}} - \sqrt{M}\mu$$

which by CLT is approximately Gaussian for large M . We can then make use of the Cramer-Slutsky Theorem, which states that if $(X_n)_{n \geq 1}$ and $(Y_n)_{n \geq 1}$ are two sequences such that $X_n \xrightarrow{d} X$ and $Y_n \xrightarrow{p} a$ as $n \rightarrow \infty$ where a is a constant, then as $n \rightarrow \infty$, it holds that $X_n * Y_n \xrightarrow{d} X * a$. Thus, Term B is approximately 0 for large M .

Term C

We have

$$\text{Term C} = \mathcal{O}\left[\sqrt{M} \left(\frac{\sum_{m=1}^M (\mathbf{W}^{(j)} \mathbf{x}_m - \mu_B)^2}{M} - \sigma^2\right)^2\right]$$

Since $E[(\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^2] = \sigma^2$ we can make the same use of Cramer-Slutsky as for Term B, such that Term C is approximately 0 for large M .

Finalizing the distribution

We have approximately

$$\sqrt{M}(\sigma_B - \sigma) \sim \mathcal{N}\left(0, \frac{E[(\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^4] - \sigma^4}{4\sigma^2}\right)$$

so

$$\sigma_B \sim \mathcal{N}\left(\sigma, \frac{E[(\mathbf{W}^{(j)} \mathbf{x}_m - \mu)^4] - \sigma^4}{4\sigma^2 M}\right)$$

1.4. Prior

Here we make use of the stochasticity from BN modeled in the Appendix section 1.3, to evaluate the implied prior on the stochastic variables for a BN network. Specifically, we consider a BN network with fully connected layers and BN applied to each layer, trained with L2-regularization (weight decay). In the following, we make use of the simplifying assumptions of no scale and shift transformations, BN applied to each layer, and independent input units to each layer.

Equivalence between the objectives of Eq. (1) and (2) requires:

$$\begin{aligned} \frac{\partial}{\partial \theta_k} \text{KL}(q_{\theta}(\omega) || p(\omega)) &= N\tau \frac{\partial}{\partial \theta_k} \Omega(\theta) \\ &= N\tau \frac{\partial}{\partial \theta_k} \sum_{l=1}^L \lambda_l \|\mathbf{W}^l\|^2 \end{aligned} \quad (5)$$

where $\theta_k \in \theta$, and θ is the set of weights in the network. To proceed with the LHS of Eq. (5) we first need to find the approximate posterior $q_{\theta}(\omega)$ that batch normalization induces. As shown in Appendix 1.3, with some weak assumptions and approximations the Central Limit Theorem (CLT) yields Gaussian distributions of the stochastic variables $\mu_{\mathbf{B}}^u, \sigma_{\mathbf{B}}^u$, for large enough M . For any BN unit u :

$$\begin{aligned} \mu_{\mathbf{B}}^u &\approx \mathcal{N}\left(\mu^u, \frac{(\sigma^u)^2}{M}\right), \\ \sigma_{\mathbf{B}}^u &\approx \mathcal{N}\left(\sigma^u, \frac{\mathbb{E}[(\mathbf{W}^{(u)}\mathbf{x} - \mu^u)^4] - (\sigma^u)^4}{4(\sigma^u)^2 M}\right) \end{aligned} \quad (6)$$

where μ^u and σ^u are *population-level* moments (i.e. moments over \mathbf{D}).

We assume that $q_{\theta}(\omega)$ and $p(\omega)$ factorize over all stochastic variables.¹⁵ We use i as an index of the set of stochastic variables. As shown in Eq. (3) in Appendix 1.2, the factorized distributions yield:

$$\text{KL}(q_{\theta}(\omega) || p(\omega)) = \sum_i \text{KL}(q_{\theta}(\omega_i) || p(\omega_i))$$

Note that each BN unit produces two $\text{KL}(q_{\theta}(\omega_i) || p(\omega_i))$ terms: one for $\omega_i = \mu_{\mathbf{B}}^u$ and one for $\omega_i = \sigma_{\mathbf{B}}^u$. We consider these terms for one particular BN unit u , and drop the index i for brevity. We use a Gaussian prior $p(\omega_i) = \mathcal{N}(\mu_p, \sigma_p^2)$ and, for consistency, use the notation $q_{\theta}(\omega_i) = \mathcal{N}(\mu_q, \sigma_q^2)$. As shown in Eq. (4) in Appendix 1.2:

$$\text{KL}(q_{\theta}(\omega_i) || p(\omega_i)) = \ln \frac{\sigma_p}{\sigma_q} + \frac{\sigma_q^2 + (\mu_q - \mu_p)^2}{2\sigma_p^2} - \frac{1}{2}$$

Since θ_k changes during training, a prior cannot depend on θ_k so $\frac{\partial}{\partial \theta_k}(\mu_p) = \frac{\partial}{\partial \theta_k}(\sigma_p) = 0$. Letting $(\cdot)'$ denote $\frac{\partial}{\partial \theta_k}(\cdot)$:

$$\frac{\partial}{\partial \theta_k} \text{KL}(q_{\theta}(\omega_i) || p(\omega_i)) = \frac{\sigma_q \sigma_q' + \mu_q \mu_q' - \mu_p \mu_q'}{\sigma_p^2} - \frac{\sigma_q'}{\sigma_q} \quad (7)$$

We need not consider θ_k past a previous layer's BN, since a normalization step is performed before scale and shift. In the general case with a given Gaussian $p(\omega)$, Eq. 7 evaluated on all BN units' means and standard deviations w.r.t. all θ_k up to a previous layer's BN, would yield an expression for a custom $N\tau \frac{\partial}{\partial \theta_k} \Omega(\theta)$ that could be used for an exact VI treatment of BN.

In our reconciliation of weight decay however, given our assumptions of no scale and shift and BN applied to each layer, we need only consider the *weights* in the same layer as the BN unit. This means that the stochastic variables in layer l are only affected by weights in $\theta_k \in \mathbf{W}^l$ (i.e. not the scale and shift variables operating on the input to the layer). We denote a weight connecting the k :th input unit to the u :th BN unit by $\mathbf{W}^{(u,k)}$. For such weights, we need to derive μ_q' and σ_q' , for two cases: $\omega_i = \mu_{\mathbf{B}}^u$ and $\omega_i = \sigma_{\mathbf{B}}^u$. We denote the priors of the mean and std. dev for $\mu_{\mathbf{B}}^u$ by $\mu_{\mu,q}$ and $\sigma_{\mu,q}$, and for $\sigma_{\mathbf{B}}^u$ by $\mu_{\sigma,q}$ and $\sigma_{\sigma,q}$. Using the distributions modeled in Eq. 6:

¹⁵The empirical distributions have been numerically checked to be linearly independent and the joint distribution is close to a bi-variate Gaussian.

Case 1: $\omega_i = \mu_{\mathbf{B}}^u$

$$\begin{aligned}\mu_{\mu,q} &= \sum_{\mathbf{x} \in \mathbf{D}} \frac{\mathbf{W}^{(u)} \mathbf{x}}{N} = \mathbf{W}^{(u)} \bar{\mathbf{x}} \\ \mu'_{\mu,q} &= \sum_{\mathbf{x} \in \mathbf{D}} \frac{\mathbf{x}_k}{N} = \bar{\mathbf{x}}_k \\ \sigma_{\mu,q} &= \sqrt{\frac{(\sigma^u)^2}{M}} = \sqrt{\frac{\sum_{\mathbf{x} \in \mathbf{D}} (\mathbf{W}^{(u)} \mathbf{x} - \mu_q)^2}{NM}} \\ \sigma'_{\mu,q} &= \frac{1}{2} \sigma_q^{-1} \sum_{\mathbf{x} \in \mathbf{D}} \frac{2(\mathbf{W}^{(u)} \mathbf{x} - \mu_q)(\mathbf{x}_k - \bar{\mathbf{x}}_k)}{NM} = \sigma_q^{-1} \left(\sum_{i=1}^K \mathbf{W}^{(u,i)} \text{Cov}(x_i, x_k) \right) M^{-1}\end{aligned}$$

where there are K input units to the layer.

Case 2: $\omega_i = \sigma_{\mathbf{B}}^u$

$$\begin{aligned}\mu_{\sigma,q} &= \sqrt{\frac{\sum_{\mathbf{x} \in \mathbf{D}} (\mathbf{W}^{(u)} \mathbf{x} - \mu_q)^2}{N}} = \sigma_{\mu,q} M^{\frac{1}{2}} \\ \mu'_{\sigma,q} &= \sigma_{\mu,q}^{-1} M^{-\frac{1}{2}} \left(\sum_{i=1}^K \mathbf{W}^{(u,i)} \text{Cov}(x_i, x_k) \right) \\ \sigma_{\sigma,q} &= \frac{\mathbb{E}[(\mathbf{W}^{(u)} \mathbf{x} - \mu^u)^4] - (\sigma^u)^4}{4(\sigma^u)^2 M} \\ \sigma'_{\sigma,q} &= \frac{\mathbb{E}[(\mathbf{W}^{(u)} \mathbf{x} - \mu^u)^4]' \sigma^u - 2(\sigma^u)^4 (\sigma^u)' - 2(\sigma^u)' \mathbb{E}[(\mathbf{W}^{(u)} \mathbf{x} - \mu^u)^4]}{4(\sigma^u)^3 M}\end{aligned}$$

Combining these results with Eq. 7 we find that taking $\text{KL}(q_{\theta}(\omega_i) || p(\omega_i))$ for the mean and variance of a single BN unit u wrt the weight from input unit k :

$$\begin{aligned}& \frac{\partial}{\partial \mathbf{W}^{(u,k)}} \text{KL}(q_{\theta}(\mu_{\mathbf{B}}^u) || p(\mu_{\mathbf{B}}^u)) + \frac{\partial}{\partial \mathbf{W}^{(u,k)}} \text{KL}(q_{\theta}(\sigma_{\mathbf{B}}^u) || p(\sigma_{\mathbf{B}}^u)) \\ &= \frac{\sigma_{\mu,q} \sigma'_{\mu,q} + \mu_{\mu,q} \mu'_{\mu,q} - \mu_{\mu,p} \mu'_{\mu,q}}{\sigma_{\mu,p}^2} - \frac{\sigma'_{\mu,q}}{\sigma_{\mu,q}} \\ &+ \frac{\sigma_{\sigma,q} \sigma'_{\sigma,q} + \mu_{\sigma,q} \mu'_{\sigma,q} - \mu_{\sigma,p} \mu'_{\sigma,q}}{\sigma_{\sigma,p}^2} - \frac{\sigma'_{\sigma,q}}{\sigma_{\sigma,q}} \\ &= \frac{\mathcal{O}(M^{-1}) + \bar{\mathbf{x}}_k \mathbf{W}^{(u)} \bar{\mathbf{x}} - \mu_{\mu,p} \bar{\mathbf{x}}_k}{\sigma_{\mu,p}^2} - \mathcal{O}(M^{-1}) \\ &+ \frac{\mathcal{O}(M^{-2}) + \sum_{i=1}^K \mathbf{W}^{(u,i)} \text{Cov}(x_i, x_k) - \mu_{\sigma,p} \mathcal{O}(M^{-\frac{1}{2}})}{\sigma_{\sigma,p}^2} \\ &- \frac{\mathbb{E}[(\mathbf{W}^{(u)} \mathbf{x} - \mu^u)^4]' \sigma^u - 2(\sigma^u)^4 (\sigma^u)' - 2(\sigma^u)' \mathbb{E}[(\mathbf{W}^{(u)} \mathbf{x} - \mu^u)^4]}{\mathbb{E}[(\mathbf{W}^{(u)} \mathbf{x} - \mu^u)^4] \sigma^u - (\sigma^u)^5}\end{aligned}$$

where we summarize the terms scaled by M with \mathcal{O} -notation. We see that if we let $M \rightarrow \infty$, $\mu_{\mu,p} = 0$, $\sigma_{\mu,p} \rightarrow \infty$, $\mu_{\sigma,p} = 0$ and $\sigma_{\sigma,p}$ is small enough, then:

$$\frac{\partial}{\partial \mathbf{W}^{(u,k)}} \left(\text{KL}(q_{\theta}(\mu_{\mathbf{B}}^u) || p(\mu_{\mathbf{B}}^u)) + \text{KL}(q_{\theta}(\sigma_{\mathbf{B}}^u) || p(\sigma_{\mathbf{B}}^u)) \right) \approx \frac{\sum_{i=1}^K \mathbf{W}^{(u,i)} \text{Cov}(x_i, x_k)}{\sigma_{\sigma,p}^2}$$

such that each BN layer yields the following:

$$\sum_u \sum_{i=1}^K \frac{\partial}{\partial \mathbf{W}^{(u,i)}} \left(\text{KL}(q_{\theta}(\mu_{\mathbf{B}}^u) || p(\mu_{\mathbf{B}}^u)) + \text{KL}(q_{\theta}(\sigma_{\mathbf{B}}^u) || p(\sigma_{\mathbf{B}}^u)) \right) \approx \sum_u \frac{\sum_{i=1}^K \mathbf{W}^{u,i} \sum_{i2=1}^K \text{Cov}(x_i, x_{i2})}{\sigma_{\sigma,p,u}^2} \quad (8)$$

where we denote the prior for the std. dev. of the std. dev. of BN unit u by $\sigma_{\sigma,p,u}$. Given our assumptions of no scale and shift from the previous layer, and independent input features in every layer, Eq. 8 reduces to:

$$\sum_u \sum_{i=1}^K \frac{\mathbf{W}^{u,i}}{\sigma_{\sigma,p}^2}$$

if the same prior is chosen for each BN unit in the layer. We therefore find that Eq. 5 is reconciled by $p(\mu_{\mathbf{B}}^u) \rightarrow \mathcal{N}(0, \infty)$ and $p(\sigma_{\mathbf{B}}^u) \rightarrow \mathcal{N}(0, \frac{1}{2N\tau\lambda_l})$, if $\frac{1}{2N\tau\lambda_l}$ is small enough, which is the case if N is large.

1.5. predictive distribution properties

This section provides derivations of properties of the predictive distribution $p^*(\mathbf{y}|\mathbf{x}, \mathbf{D})$ in section 3.4, following (Gal, 2016). We first find the first two modes of the approximate predictive distribution (with the second mode applicable to regression), then show how to estimate the predictive log likelihood, a measure of uncertainty quality used in the evaluation.

Predictive mean Assuming Gaussian iid noise defined by model precision τ , i.e. $f_{\omega}(\mathbf{x}, \mathbf{y}) = p(\mathbf{y}|f_{\omega}(\mathbf{x})) = \mathcal{N}(\mathbf{y}; f_{\omega}(\mathbf{x}), \tau^{-1}\mathbf{I})$:

$$\begin{aligned} \mathbb{E}_{p^*}[\mathbf{y}] &= \int \mathbf{y} p^*(\mathbf{y}|\mathbf{x}, \mathbf{D}) d\mathbf{y} \\ &= \int_{\mathbf{y}} \mathbf{y} \left(\int_{\omega} f_{\omega}(\mathbf{x}, \mathbf{y}) q_{\theta}(\omega) d\omega \right) d\mathbf{y} \\ &= \int_{\mathbf{y}} \mathbf{y} \left(\int_{\omega} \mathcal{N}(\mathbf{y}; f_{\omega}(\mathbf{x}), \tau^{-1}\mathbf{I}) q_{\theta}(\omega) d\omega \right) d\mathbf{y} \\ &= \int_{\omega} \left(\int_{\mathbf{y}} \mathbf{y} \mathcal{N}(\mathbf{y}; f_{\omega}(\mathbf{x}), \tau^{-1}\mathbf{I}) d\mathbf{y} \right) q_{\theta}(\omega) d\omega \\ &= \int_{\omega} f_{\omega}(\mathbf{x}) q_{\theta}(\omega) d\omega \\ &\approx \frac{1}{T} \sum_{i=1}^T f_{\hat{\omega}_i}(\mathbf{x}) \end{aligned}$$

where we take the MC Integral with T samples of ω for the approximation in the final step.

Predictive variance For regression, our goal is to estimate:

$$\text{Cov}_{p^*}[\mathbf{y}] = \mathbb{E}_{p^*}[\mathbf{y}^T \mathbf{y}] - \mathbb{E}_{p^*}[\mathbf{y}]^T \mathbb{E}_{p^*}[\mathbf{y}]$$

We find that:

$$\begin{aligned} \mathbb{E}_{p^*}[\mathbf{y}^T \mathbf{y}] &= \int_{\mathbf{y}} \mathbf{y}^T \mathbf{y} p^*(\mathbf{y}|\mathbf{x}, \mathbf{D}) d\mathbf{y} \\ &= \int_{\mathbf{y}} \mathbf{y}^T \mathbf{y} \left(\int_{\omega} f_{\omega}(\mathbf{x}, \mathbf{y}) q_{\theta}(\omega) d\omega \right) d\mathbf{y} \\ &= \int_{\omega} \left(\int_{\mathbf{y}} \mathbf{y}^T \mathbf{y} f_{\omega}(\mathbf{x}, \mathbf{y}) d\mathbf{y} \right) q_{\theta}(\omega) d\omega \\ &= \int_{\omega} \left(\text{Cov}_{f_{\omega}(\mathbf{x}, \mathbf{y})}(\mathbf{y}) + \mathbb{E}_{f_{\omega}(\mathbf{x}, \mathbf{y})}[\mathbf{y}]^T \mathbb{E}_{f_{\omega}(\mathbf{x}, \mathbf{y})}[\mathbf{y}] \right) q_{\theta}(\omega) d\omega \\ &= \int_{\omega} \left(\tau^{-1}\mathbf{I} + f_{\omega}(\mathbf{x})^T f_{\omega}(\mathbf{x}) \right) q_{\theta}(\omega) d\omega \\ &= \tau^{-1}\mathbf{I} + E_{q_{\theta}(\omega)}[f_{\omega}(\mathbf{x})^T f_{\omega}(\mathbf{x})] \\ &\approx \tau^{-1}\mathbf{I} + \frac{1}{T} \sum_{i=1}^T f_{\hat{\omega}_i}(\mathbf{x})^T f_{\hat{\omega}_i}(\mathbf{x}) \end{aligned}$$

where we use MC integration with T samples for the final step. The predictive covariance matrix is given by:

$$\text{Cov}_{p^*}[\mathbf{y}] \approx \tau^{-1}\mathbf{I} + \frac{1}{T} \sum_{i=1}^T f_{\hat{\omega}_i}(\mathbf{x})^\top f_{\hat{\omega}_i}(\mathbf{x}) - \mathbb{E}_{p^*}[\mathbf{y}]^\top \mathbb{E}_{p^*}[\mathbf{y}]$$

which is the sum of the variance from observation noise and the sample covariance from T stochastic forward passes through the network.

The form of p^* can be approximated by a Gaussian for each output dimension (for regression). We assume bounded domains for each input dimension, wide layers throughout the network, and a uni-modal distribution of weights centered at 0. By the Liapounov CLT condition, the first layer then receives approximately Gaussian inputs (a proof can be found in (Lehmann, 1999)). Having sampled $\mu_{\mathbf{B}}^u$ and $\sigma_{\mathbf{B}}^u$ from a mini-batch, each BN unit's output is bounded. CLT thereby continues to hold for deeper layers, including $f_{\omega}(\mathbf{x}) = \mathbf{W}^L \mathbf{x}^L$. A similar motivation for a Gaussian approximation of Dropout has been presented by (Wang & Manning, 2013).

Predictive Log Likelihood We use the Predictive Log Likelihood (PLL) as a measure to estimate the model's uncertainty quality. For a certain test point $(\mathbf{y}_i, \mathbf{x}_i)$, the PLL definition and approximation can be expressed as:

$$\begin{aligned} \text{PLL}(f_{\omega}(\mathbf{x}), (\mathbf{y}_i, \mathbf{x}_i)) &= \log p(\mathbf{y}_i | f_{\omega}(\mathbf{x}_i)) \\ &= \log \int f_{\omega}(\mathbf{x}_i, \mathbf{y}_i) p(\omega | \mathbf{D}) d\omega \\ &\approx \log \int f_{\omega}(\mathbf{x}_i, \mathbf{y}_i) q_{\theta}(\omega) d\omega \\ &\approx \log \frac{1}{T} \sum_{j=1}^T p(\mathbf{y}_i | f_{\hat{\omega}_j}(\mathbf{x}_i)) \end{aligned}$$

where $\hat{\omega}_j$ represents a sampled set of stochastic parameters from the approximate posterior distribution $q_{\theta}(\omega)$ and we take a MC integration with T samples. For regression, due to the iid Gaussian noise, we can further develop the derivation into the form we use when sampling:

$$\begin{aligned} \text{PLL}(f_{\omega}(\mathbf{x}), (\mathbf{y}_i, \mathbf{x}_i)) &= \log \frac{1}{T} \sum_{j=1}^T \mathcal{N}(\mathbf{y}_i | f_{\hat{\omega}_j}(\mathbf{x}_i), \tau^{-1}\mathbf{I}) \\ &= \log \text{sumexp}_{j=1, \dots, T} \left(-\frac{1}{2} \tau \| \mathbf{y}_i - f_{\hat{\omega}_j}(\mathbf{x}_i) \|^2 \right) \\ &\quad - \log T - \frac{1}{2} \log 2\pi + \frac{1}{2} \log \tau \end{aligned}$$

Note that PLL makes no assumption on the form of the approximate predictive distribution.

1.6. Data

To assess the uncertainty quality of the various methods studied we rely on eight standard regression datasets, listed in Table 3. Publicly available from the UCI Machine Learning Repository (University of California, 2017) and Delve (Ghahramani, 1996), these datasets have been used to benchmark comparative models in recent related literature (see (Hernández-Lobato & Adams, 2015), (Gal & Ghahramani, 2015), (Bui et al., 2016) and (Li & Gal, 2017)).

For image classification, we applied MCBN using ResNet32 to CIFAR10.

For the image segmentation task, we applied MCBN using Bayesian SegNet on data from CamVid and PASCAL-VOC using models published in (Kendall et al., 2015).

Table 3. **Regression dataset summary.** Properties of the eight regression datasets used to evaluate MCBN. N is the dataset size and Q is the n.o. input features. Only one target feature was used – we used heating load for the Energy Efficiency dataset, which contains multiple target features.

Dataset name	N	Q
Boston Housing	506	13
Concrete Compressive Strength	1,030	8
Energy Efficiency	768	8
Kinematics 8nm	8,192	8
Power Plant	9,568	4
Protein Tertiary Structure	45,730	9
Wine Quality (Red)	1,599	11
Yacht Hydrodynamics	308	6

1.7. Extended experimental results

Below, we provide extended results measuring uncertainty quality. In Tables 4 and 5, we provide tables showing the mean $\overline{\text{CRPS}}$ and $\overline{\text{PLL}}$ values for MCBN and MCDO. These results indicate that MCBN performs on par or better than MCDO across several datasets. In Table 6 we provide the raw PLL and CRPS results for MCBN and MCDO. In Table 7 we provide RMSE results of the MCBN and MCDO networks in comparison with non-stochastic BN and DO networks. These results indicate that the procedure of multiple forward passes in MCBN and MCDO show slight improvements in the accuracy of the network.

In Figure 4 and Figure 5, we provide a full set of our uncertainty quality visualization plots, where errors in predictions are sorted by estimated uncertainty. The shaded areas show the model uncertainty and gray dots show absolute prediction errors on the test set. A gray line depicts a running mean of the errors. The dashed line indicates the optimized constant uncertainty. In these plots, we can see a correlation between estimated uncertainty (shaded area) and mean error (gray). This trend indicates that the model uncertainty estimates can recognize samples with larger (or smaller) potential for predictive errors.

We also conduct a sensitivity analysis to estimate how the uncertainty quality varies with batch size M and the number of stochastic forward passes T . In tables 8 and 9 we evaluate $\overline{\text{CRPS}}$ and $\overline{\text{PLL}}$ respectively for the regression datasets when trained and evaluated with varying batch sizes, but other hyperparameters fixed (T was fixed at 100). The results show that results deteriorate when batch sizes are too small, likely stemming from the large variance of the approximate posterior. In tables 10 and 11 we evaluate $\overline{\text{CRPS}}$ and $\overline{\text{PLL}}$ respectively for the regression datasets when trained and evaluated with varying n.o. stochastic forward samples, but other hyperparameters fixed (M was fixed at 128). The results are indicative of performance improvements with larger T , although we see improvements over baseline for some datasets already with $T = 50$ (1/10:th of the T used in our main experiments).

Table 4. **Uncertainty quality measured by $\overline{\text{CRPS}}$ on regression datasets.** $\overline{\text{CRPS}}$ measured on eight datasets over 5 random 80-20 splits of the data with 5 different random seeds each split. Mean values for MCBN, MCDO and MNF are reported along with standard error. A significance test was performed to check if CRPS significantly exceeds the baseline. The p -value from a one sample t-test is reported.

Dataset	$\overline{\text{CRPS}}$					
	MCBN	p -value	MCDO	p -value	MNF	p -value
Boston Housing	8.50±0.86	6.39E-10	3.06±0.33	1.64E-09	5.88±1.09	2.01E-05
Concrete	3.91±0.25	4.53E-14	0.93±0.41	3.13E-02	3.13±0.81	6.43E-04
Energy Efficiency	5.75±0.52	6.71E-11	1.37±0.89	1.38E-01	1.10±2.63	6.45E-01
Kinematics 8nm	2.85±0.18	2.33E-14	1.82±0.14	1.64E-12	0.52±0.26	7.15E-02
Power Plant	0.24±0.05	2.32E-04	-0.44±0.05	2.17E-08	-0.89±0.15	3.36E-06
Protein	2.66±0.10	2.77E-12	0.99±0.08	2.34E-12	0.57±0.03	8.56E-16
Wine Quality (Red)	0.26±0.07	1.26E-03	2.00±0.21	1.83E-09	0.93±0.12	6.19E-08
Yacht Hydrodynamics	-56.39±14.27	5.94E-04	21.42±2.99	2.16E-07	24.92±3.77	9.62E-06

Bayesian Uncertainty Estimation for Batch Normalized Deep Networks

Table 5. Uncertainty quality measured by $\overline{\text{PLL}}$ on regression datasets. $\overline{\text{PLL}}$ measured on eight datasets over 5 random 80-20 splits of the data with 5 different random seeds each split. Mean values for MCBN, MCDO and MNF are reported along with standard error. A significance test was performed to check if $\overline{\text{PLL}}$ significantly exceeds the baseline. The p -value from a one sample t-test is reported.

Dataset	MCBN	p -value	$\overline{\text{PLL}}$		MNF	p -value
			MCDO	p -value		
Boston Housing	10.49±1.35	5.41E-08	5.51±1.05	2.20E-05	1.76±1.12	1.70E-01
Concrete	-36.36±12.12	6.19E-03	10.92±1.78	2.34E-06	-2.16±4.19	6.79E-01
Energy Efficiency	10.89±1.16	1.79E-09	-14.28±5.15	1.06E-02	-33.88±29.57	2.70E-01
Kinematics 8nm	1.68±0.37	1.29E-04	-0.26±0.18	1.53E-01	0.42±0.43	2.70E-01
Power Plant	0.33±0.14	2.72E-02	3.52±0.23	1.12E-13	-0.86±0.15	7.33E-06
Protein	2.56±0.23	4.28E-11	6.23±0.19	2.57E-21	0.52±0.07	1.81E-07
Wine Quality (Red)	0.19±0.09	3.72E-02	2.91±0.35	1.84E-08	0.83±0.16	2.27E-05
Yacht Hydrodynamics	45.58±5.18	5.67E-09	-41.54±31.37	1.97E-01	46.19±4.45	2.47E-07

Table 6. Raw (unnormalized) CRPS and PLL scores on regression datasets. CRPS and PLL measured on eight datasets over 5 random 80-20 splits of the data with 5 different random seeds each split. Mean values and standard errors are reported for MCBN, MCDO and MNF.

Dataset	CRPS			PLL		
	MCBN	MCDO	MNF	MCBN	MCDO	MNF
Boston Housing	1.45±0.02	1.41±0.02	1.57±0.02	-2.38±0.02	-2.35±0.02	-2.51±0.06
Concrete	2.40±0.04	2.42±0.04	3.61±0.02	-3.45±0.11	-2.94±0.02	-3.35±0.04
Energy Efficiency	0.33±0.01	0.26±0.00	1.33±0.04	-0.94±0.04	-0.80±0.04	-3.18±0.07
Kinematics 8nm	0.04±0.00	0.04±0.00	0.05±0.00	1.21±0.01	1.24±0.00	1.04±0.00
Power Plant	2.00±0.01	2.00±0.01	2.31±0.01	-2.75±0.00	-2.72±0.01	-2.86±0.01
Protein	1.95±0.01	1.95±0.00	2.25±0.01	-2.73±0.00	-2.70±0.00	-2.83±0.01
Wine Quality (Red)	0.34±0.00	0.33±0.00	0.34±0.00	-0.95±0.01	-0.89±0.01	-0.93±0.00
Yacht Hydrodynamics	0.68±0.02	0.32±0.01	0.94±0.01	-1.39±0.03	-2.57±0.69	-1.96±0.05

Table 7. Prediction accuracy measured by RMSE on regression datasets. RMSE measured on eight datasets over 5 random 80-20 splits of the data with 5 different random seeds each split. Mean values and standard errors are reported for for MCBN, MCDO and MNF as well as conventional non-Bayesian models BN and DO.

Dataset	MCBN	BN	RMSE		
			MCDO	DO	MNF
Boston Housing	2.75±0.05	2.77±0.05	2.65±0.05	2.69±0.05	2.98±0.06
Concrete	4.78±0.09	4.89±0.08	4.80±0.10	4.99±0.10	6.57±0.04
Energy Efficiency	0.59±0.02	0.57±0.01	0.47±0.01	0.49±0.01	2.38±0.07
Kinematics 8nm	0.07±0.00	0.07±0.00	0.07±0.00	0.07±0.00	0.09±0.00
Power Plant	3.74±0.01	3.74±0.01	3.74±0.02	3.72±0.02	4.19±0.01
Protein	3.66±0.01	3.69±0.01	3.66±0.01	3.68±0.01	4.10±0.01
Wine Quality (Red)	0.62±0.00	0.62±0.00	0.60±0.00	0.61±0.00	0.61±0.00
Yacht Hydrodynamics	1.23±0.05	1.28±0.06	0.75±0.03	0.72±0.04	2.13±0.05

Bayesian Uncertainty Estimation for Batch Normalized Deep Networks

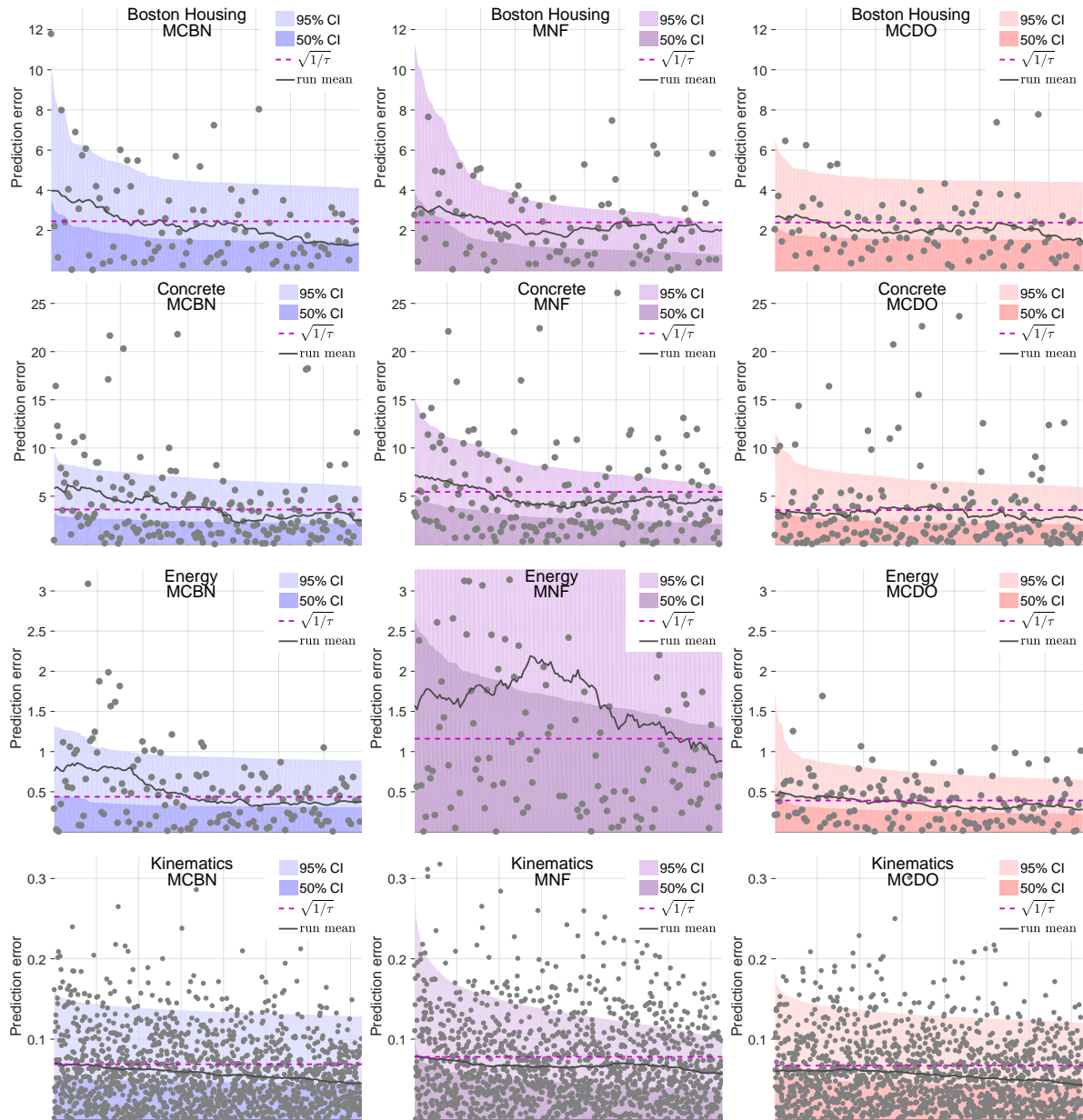


Figure 4. Errors in predictions (gray dots) sorted by estimated uncertainty on select datasets. The shaded areas show model uncertainty (light area 95% CI, dark area 50% CI). Gray dots show absolute prediction errors on the test set, and the gray line depicts a running mean of the errors. The dashed line indicates the optimized constant uncertainty. A correlation between estimated uncertainty (shaded area) and mean error (gray) indicates the uncertainty estimates are meaningful for estimating errors.

Bayesian Uncertainty Estimation for Batch Normalized Deep Networks

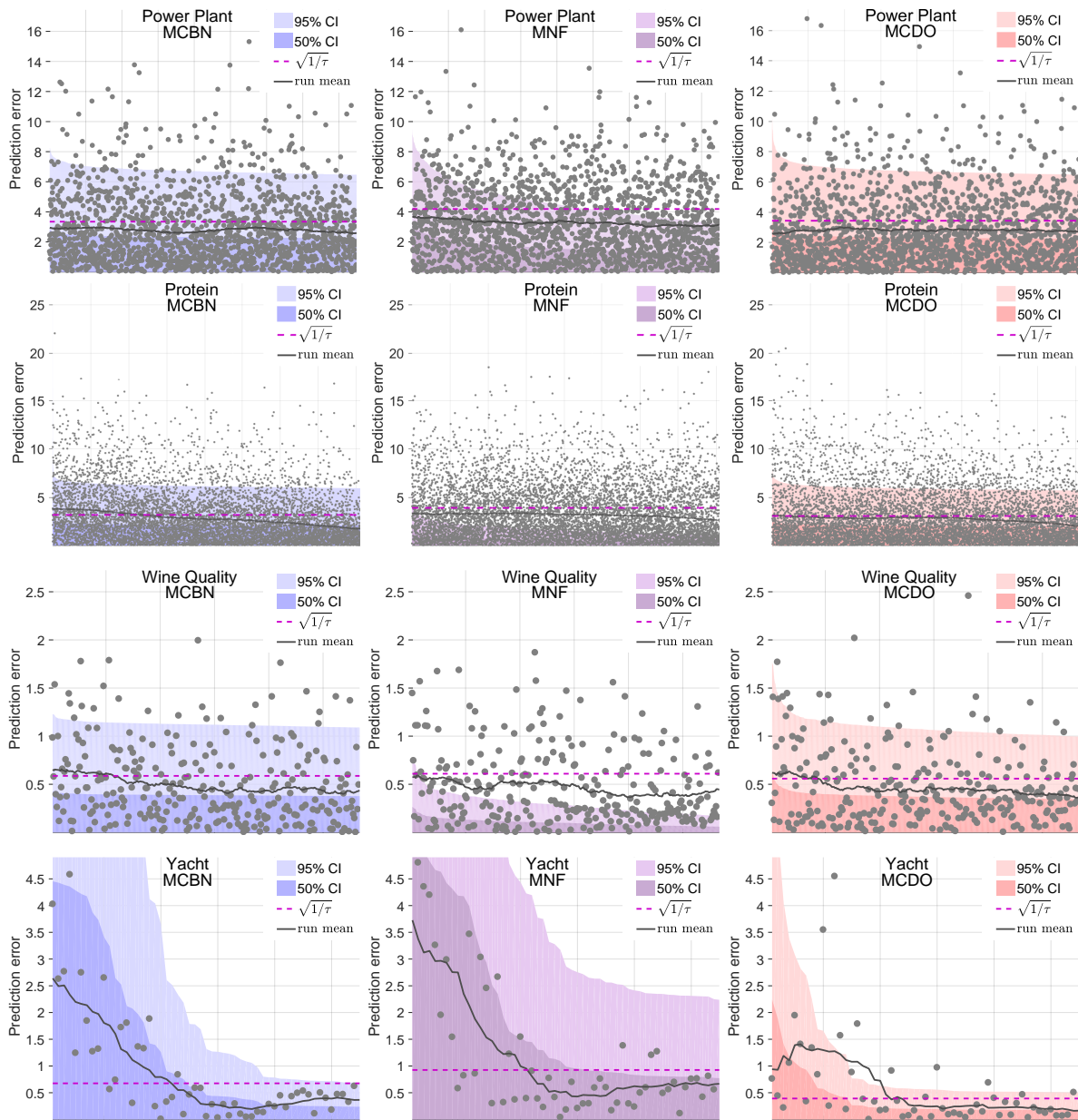


Figure 5. Errors in predictions (gray dots) sorted by estimated uncertainty on select datasets. The shaded areas show model uncertainty (light area 95% CI, dark area 50% CI). Gray dots show absolute prediction errors on the test set, and the gray line depicts a running mean of the errors. The dashed line indicates the optimized constant uncertainty. A correlation between estimated uncertainty (shaded area) and mean error (gray) indicates the uncertainty estimates are meaningful for estimating errors.

Table 8. **Uncertainty quality sensitivity to batch size.** A sensitivity analysis to determine how MCBN uncertainty quality varies with batch size is measured on eight regression datasets using $\overline{\text{CRPS}}$ as the quality measure. Results are measured over 3 random 80-20 splits of the data with 5 different random seeds each split.

Batch size	$\overline{\text{CRPS}}$							
	8	16	32	64	128	256	512	1024
Boston Housing	-7.1	16.6	11.8	7.2	2.5	0.9	-	-
Concrete	-34.5	6.0	5.0	5.1	2.9	1.4	0.6	0.0
Energy Efficiency	-61.6	-3.0	2.7	9.8	11.1	0.8	4.9	-
Kinematics 8nm	-1.4	-4.3	0.2	2.8	2.7	1.7	0.9	0.5
Power Plant	-10.5	0.8	0.0	-0.1	0.0	0.0	0.2	0.0
Protein	14.5	4.8	3.6	2.8	2.5	1.6	1.0	0.5
Wine Quality (Red)	2.2	1.6	0.6	0.6	0.3	0.0	0.2	0.0
Yacht Hydrodynamics	15.1	-23.0	-30.4	21.0	34.4	-	-	-

Table 9. **Uncertainty quality sensitivity to batch size.** A sensitivity analysis to determine how MCBN uncertainty quality varies with batch size is measured on eight regression datasets using $\overline{\text{PLL}}$ as the quality measure. Results are measured over 3 random 80-20 splits of the data with 5 different random seeds each split.

Batch size	$\overline{\text{PLL}}$							
	8	16	32	64	128	256	512	1024
Boston Housing	13.9	-36.7	10.0	7.9	3.7	1.5	-	-
Concrete	-113.3	-528.4	-10.0	2.9	0.0	1.4	0.2	0.0
Energy Efficiency	-64.4	5.2	-0.2	-9.6	-14.5	1.4	10.4	-
Kinematics 8nm	-4.9	-5.4	-3.1	1.6	2.3	1.5	0.7	0.4
Power Plant	-135.0	-1.4	-1.0	-1.1	-0.4	0.1	-0.1	0.4
Protein	44.9	15.7	4.6	2.9	2.8	2.2	1.2	0.6
Wine Quality (Red)	2.2	2.0	0.0	0.5	0.6	0.4	0.0	0.0
Yacht Hydrodynamics	99.6	74.9	76.8	48.5	44.9	-	-	-

Table 10. **Uncertainty quality sensitivity to n.o. stochastic forward passes.** A sensitivity analysis to determine how MCBN uncertainty quality varies with the n.o. stochastic forward passes measured on eight regression datasets using $\overline{\text{CRPS}}$ as the quality measure. Results are measured over 3 random 80-20 splits of the data with 5 different random seeds each split.

Forward passes	$\overline{\text{CRPS}}$		
	250	100	50
Boston Housing	6.1	2.7	3.2
Concrete	3.3	2.3	3.3
Energy Efficiency	13.2	4.2	7.9
Kinematics 8nm	3.2	2.7	4.2
Power Plant	0.2	0.5	0.1
Protein	2.3	2.7	2.4
Wine Quality (Red)	0.9	-0.4	0.6
Yacht Hydrodynamics	32.9	32.2	32.1

Table 11. Uncertainty quality sensitivity to n.o. stochastic forward passes. A sensitivity analysis to determine how MCBN uncertainty quality varies with the n.o. stochastic forward passes measured on eight regression datasets using $\overline{\text{PLL}}$ as the quality measure. Results are measured over 3 random 80-20 splits of the data with 5 different random seeds each split.

Forward passes	$\overline{\text{PLL}}$		
	250	100	50
Boston Housing	7.8	1.9	2.6
Concrete	3.8	7.1	0.1
Energy Efficiency	15.7	-30.5	-47.3
Kinematics 8nm	2.5	2.2	3.4
Power Plant	-0.9	0.7	-0.9
Protein	1.8	2.0	2.4
Wine Quality (Red)	1.7	-0.9	1.1
Yacht Hydrodynamics	38.0	35.9	35.5

1.8. Uncertainty in image segmentation

We applied MCBN to an image segmentation task using Bayesian SegNet with the main CamVid and PASCAL-VOC models in (Kendall et al., 2015). Here, we provide more image from Pascal VOC dataset in Figure 6.

1.9. Batch normalization statistics

In Figure 7 and Figure 8, we provide statistics on the batch normalization parameters used for training. The plots show the distribution of BN mean and BN variance over different mini-batches of an actual training of Yacht dataset for one unit in the first hidden layer and the second hidden layer. Data is provided for different epochs and for different batch sizes.

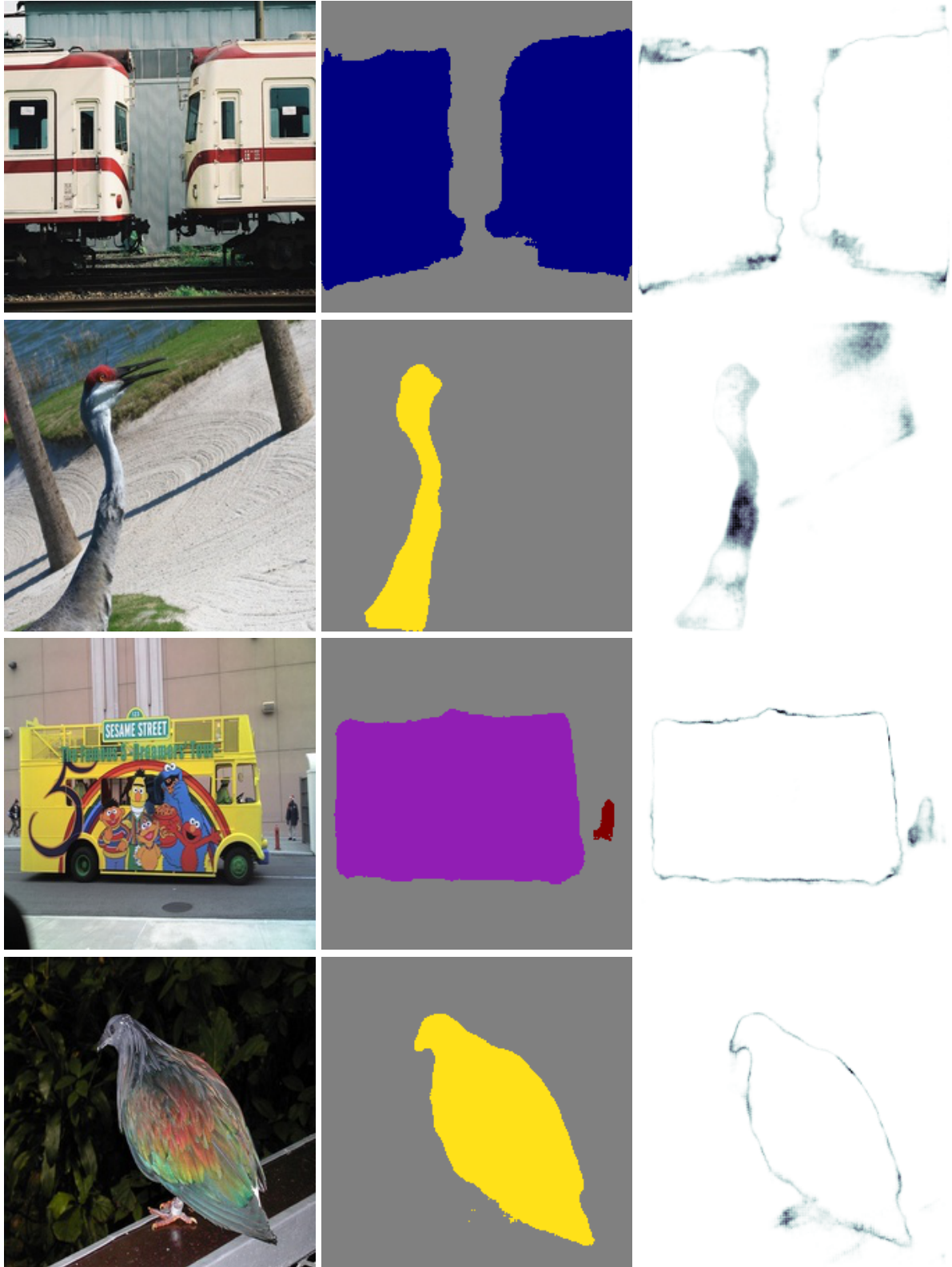


Figure 6. **Uncertainty in image segmentation.** Results applying MCBN to Bayesian SegNet (Kendall et al., 2015) on images from PASCAL-VOC (right). Left: original. Middle: the Bayesian estimated segmentation. Right: estimated uncertainty using MCBN for all classes. Mini-batches of size 36 were used for PASCAL-VOC on images of size 224x224. 20 inferences were conducted to estimate the mean and variance of MCBN.

Bayesian Uncertainty Estimation for Batch Normalized Deep Networks

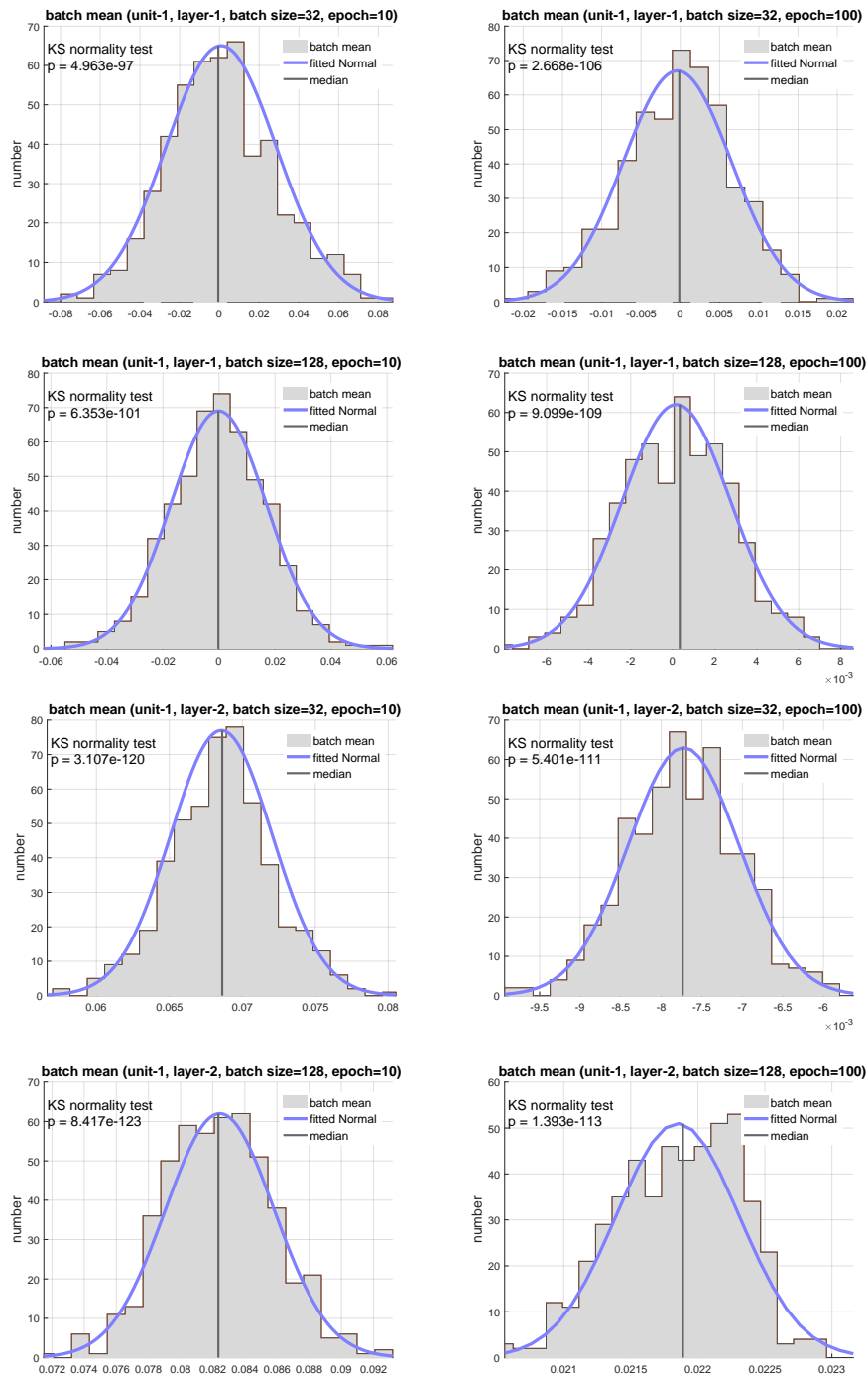


Figure 7. The distribution of means of mini-batches during training of one of our datasets. The distribution closely follows our analytically approximated Gaussian distribution. The data is collected for one unit of each layer and is provided for different epochs and for different batch sizes.

Bayesian Uncertainty Estimation for Batch Normalized Deep Networks

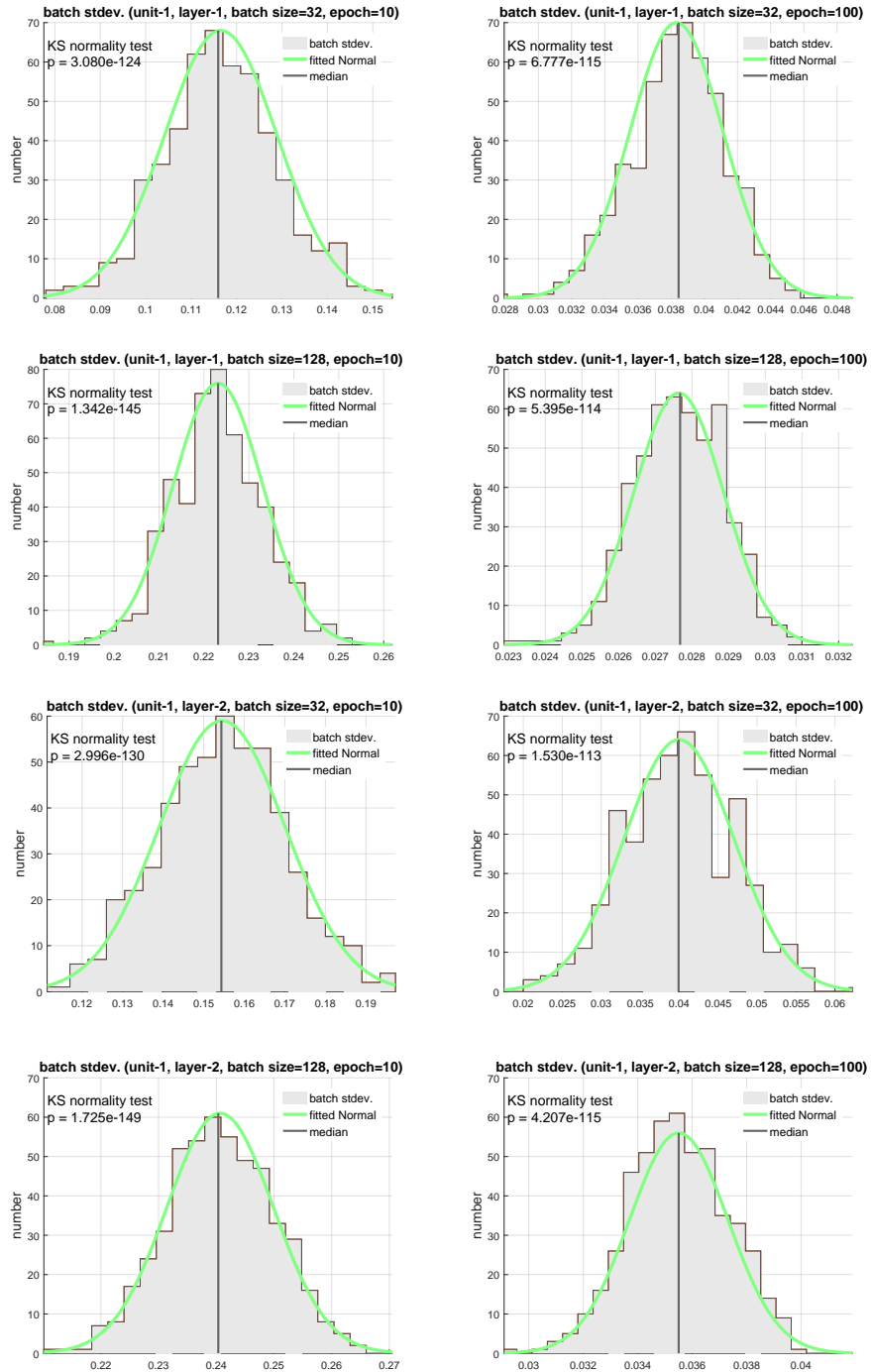


Figure 8. The distribution of standard deviation of mini-batches during training of one of our datasets. The distribution closely follows our analytically approximated Gaussian distribution. The data is collected for one unit of each layer and is provided for different epochs and for different batch sizes.

References

- Bui, T. D., Hernández-Lobato, D., Li, Y., Hernández-Lobato, J. M., and Turner, R. E. Deep Gaussian Processes for Regression using Approximate Expectation Propagation. In *ICML*, 2016.
- Chen, X., Kundu, K., Zhang, Z., Ma, H., Fidler, S., and Urtasun, R. Monocular 3d object detection for autonomous driving. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2147–2156, 2016.
- Djuric, U., Zadeh, G., Aldape, K., and Diamandis, P. Precision histology: how deep learning is poised to revitalize histomorphology for personalized cancer care. *npj Precision Oncology*, 1(1):22, 2017.
- Esteva, A., Kuprel, B., Novoa, R. A., Ko, J., Swetter, S. M., Blau, H. M., and Thrun, S. Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, Feb 2017.
- Gal, Y. *Uncertainty in Deep Learning*. PhD thesis, University of Cambridge, 2016.
- Gal, Y. and Ghahramani, Z. Dropout as a Bayesian Approximation : Representing Model Uncertainty in Deep Learning. *ICML*, 48:1–10, 2015.
- Ghahramani, Z. *Delve Datasets*. University of Toronto, 1996. URL <http://www.cs.toronto.edu/{~}delve/data/kin/desc.html>.
- Ghahramani, Z. Probabilistic machine learning and artificial intelligence. *Nature*, 521(7553):452–459, May 2015.
- Gneiting, T. and Raftery, A. E. Strictly Proper Scoring Rules, Prediction, and Estimation. *Journal of the American Statistical Association*, 102(477):359–378, 2007.
- Goodfellow, I. J., Shlens, J., and Szegedy, C. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- Graves, A. Practical Variational Inference for Neural Networks. *NIPS*, 2011.
- Hernández-Lobato, J. M. and Adams, R. Probabilistic backpropagation for scalable learning of bayesian neural networks. In *International Conference on Machine Learning*, pp. 1861–1869, 2015.
- Hinton, G. E. and Van Camp, D. Keeping the neural networks simple by minimizing the description length of the weights. In *Proceedings of the sixth annual conference on Computational learning theory*, pp. 5–13. ACM, 1993.
- Ioffe, S. Batch renormalization: Towards reducing minibatch dependence in batch-normalized models. *CoRR*, abs/1702.03275, 2017. URL <http://arxiv.org/abs/1702.03275>.
- Ioffe, S. and Szegedy, C. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. *Arxiv*, 2015. URL <http://arxiv.org/abs/1502.03167>.
- Karpathy, A. Convnetjs demo: toy 1d regression, 2015. URL <http://cs.stanford.edu/people/karpathy/convnetjs/demo/regression.html>.
- Kendall, A., Badrinarayanan, V., and Cipolla, R. Bayesian SegNet: Model Uncertainty in Deep Convolutional Encoder-Decoder Architectures for Scene Understanding. *CoRR*, abs/1511.0, 2015. URL <http://arxiv.org/abs/1511.02680>.
- Kingma, D. P. and Welling, M. Auto-Encoding Variational Bayes. In *ICLR*, 2014.
- Krizhevsky, A. and Hinton, G. Learning multiple layers of features from tiny images. 2009.
- Krueger, D., Huang, C.-W., Islam, R., Turner, R., Lacoste, A., and Courville, A. Bayesian hypernetworks. *arXiv preprint arXiv:1710.04759*, 2017.
- Lehmann, E. L. *Elements of Large-Sample Theory*. Springer Verlag, New York, 1999. ISBN 0387985956.
- Li, Y. and Gal, Y. Dropout Inference in Bayesian Neural Networks with Alpha-divergences. *arXiv*, 2017.

- Louizos, C. and Welling, M. Multiplicative normalizing flows for variational Bayesian neural networks. In Precup, D. and Teh, Y. W. (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 2218–2227, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR. URL <http://proceedings.mlr.press/v70/louizos17a.html>.
- MacKay, D. J. A practical bayesian framework for backpropagation networks. *Neural computation*, 4(3):448–472, 1992.
- Neal, R. M. *Bayesian Learning for Neural Networks*. PhD thesis, University of Toronto, 1995.
- Neal, R. M. *Bayesian learning for neural networks*, volume 118. Springer Science & Business Media, 2012.
- Selten, R. Axiomatic characterization of the quadratic scoring rule. *Experimental Economics*, 1(1):43–62, 1998.
- Shen, L. End-to-end training for whole image breast cancer diagnosis using an all convolutional design. *arXiv preprint arXiv:1708.09427*, 2017.
- University of California, I. UC Irvine Machine Learning Repository, 2017. URL <https://archive.ics.uci.edu/ml/index.html>.
- Wang, S. I. and Manning, C. D. Fast dropout training. *Proceedings of the 30th International Conference on Machine Learning*, 28:118–126, 2013. URL <http://machinelearning.wustl.edu/mlpapers/papers/wang13a>.