

---

# Training Binary Neural Networks using the Bayesian Learning Rule

---

Xiangming Meng<sup>1</sup> Roman Bachmann<sup>\*2</sup> Mohammad Emtiyaz Khan<sup>1</sup>

## Abstract

Neural networks with binary weights are computation-efficient and hardware-friendly, but their training is challenging because it involves a discrete optimization problem. Surprisingly, ignoring the discrete nature of the problem and using gradient-based methods, such as the Straight-Through Estimator, still works well in practice. This raises the question: are there principled approaches which justify such methods? In this paper, we propose such an approach using the *Bayesian learning rule*. The rule, when applied to estimate a Bernoulli distribution over the binary weights, results in an algorithm which justifies some of the algorithmic choices made by the previous approaches. The algorithm not only obtains state-of-the-art performance, but also enables uncertainty estimation for continual learning to avoid catastrophic forgetting. Our work provides a principled approach for training binary neural networks which justifies and extends existing approaches.

## 1. Introduction

Deep neural networks (DNNs) have been remarkably successful in machine learning but their training and deployment requires a high energy budget and hinders their application to resource-constrained devices, such as mobile phones, wearables, and IoT devices. Binary neural networks (BiNNs), where weights and/or activations are restricted to binary values, are one promising solution to address this issue (Courbariaux et al., 2016; 2015). Compared to full-precision DNNs, e.g., using 32-bits, using BiNNs directly gives a 32 times reduction in the model size. Further computational efficiency is obtained by using specialized hardware,

e.g., by replacing the multiplication and addition operations with the bit-wise `xnor` and `bitcount` operations (Rastegari et al., 2016; Mishra et al., 2018; Bethge et al., 2020). In the near future, BiNNs are expected to play an important role in energy-efficient and hardware-friendly deep learning.

A problem with BiNNs is that their training is much more difficult than their continuous counterpart. BiNNs obtained by quantizing already trained DNNs do not work well, and it is preferable to optimize for binary weights directly. Such training is challenging because it involves a discrete optimization problem. Continuous optimization methods such as the Adam optimizer (Kingma & Ba, 2015) are not expected to perform well or even converge.

Despite such theoretical issues, a method called Straight-Through-Estimator (STE) (Bengio et al., 2013), which employs continuous optimization methods, works remarkably well (Courbariaux et al., 2015). The method is justified based on “latent” real-valued weights which are discretized at every iteration to get binary weights. The gradients used to update the latent weights, however, are computed at the binary weights (see Figure 1 (a) for an illustration). It is not clear why these gradients help the search for the minimum of the discrete problem (Yin et al., 2019; Alizadeh et al., 2019). Another recent work by Helwegen et al. (2019) dismisses the idea of latent weights, and proposes a new optimizer called Binary Optimizer (Bop) based on *inertia*. Unfortunately, the steps used by their optimizers too are derived based on intuition and are not theoretically justified using an optimization problem. Our goal in this paper is to address this issue and propose a principled approach to justify the algorithmic choices of these previous approaches.

In this paper, we present a Bayesian perspective to justify previous approaches. Instead of optimizing a discrete objective, the Bayesian approach relaxes it by using a distribution over the binary variable, resulting in a principled approach for discrete optimization. We use a Bernoulli approximation to the posterior and estimate it using a recently proposed approximate Bayesian inference method called the Bayesian learning rule (Khan & Lin, 2017; Khan & Rue, 2020). This results in an algorithm which justifies some of the algorithmic choices made by existing methods; see Table 1 for a summary of results. Since our algorithm is based on a well-defined optimization problem, it is easier to extend its

<sup>1</sup>RIKEN Center for Advanced Intelligence Project (AIP), Tokyo, Japan. <sup>2</sup>cole polytechnique fdrale de Lausanne (EPFL), Lausanne, Switzerland. Correspondence to: Mohammad Emtiyaz Khan <emtiyaz.khan@riken.jp>.

\*This work is performed during an internship at the RIKEN Center for AIP.

Proceedings of the 37<sup>th</sup> International Conference on Machine Learning, Online, PMLR 119, 2020. Copyright 2020 by the author(s).

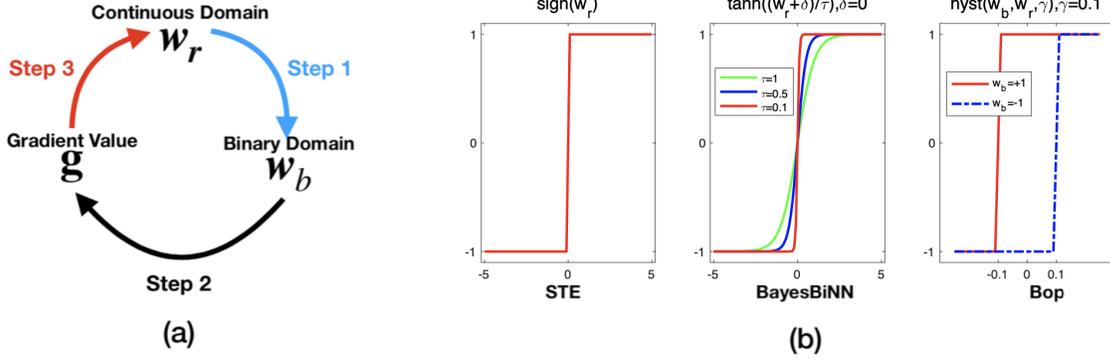


Figure 1. (a): The three steps involved in training BiNNs. In step 1, we obtain binary weights  $w_b$  from the real-valued parameters  $w_r$ . In step 2, we compute gradients at  $w_b$  and, in step 3, update  $w_r$ . (b): Different functions used to convert continuous parameters to binary weights. From left to right: Sign function used in STE; Tanh function used in BayesBiNN; Hysteresis function (see (25)) used in Bop.

	STE	Our BayesBiNN method	Bop
Step 1: Get $w_b$ from $w_r$	$w_b \leftarrow \text{sign}(w_r)$	$w_b \leftarrow \tanh((w_r + \delta)/\tau)$	$w_b \leftarrow \text{hyst}(w_r, w_b, \gamma)$
Step 2: Compute gradient at $w_b$	$\mathbf{g} \leftarrow \nabla_{w_b} \ell(y, f_{w_b}(\mathbf{x}))$	$\mathbf{g} \leftarrow \nabla_{w_b} \ell(y, f_{w_b}(\mathbf{x}))$	$\mathbf{g} \leftarrow \nabla_{w_b} \ell(y, f_{w_b}(\mathbf{x}))$
Step 3: Update $w_r$	$w_r \leftarrow w_r - \alpha \mathbf{g}$	$w_r \leftarrow (1 - \alpha)w_r - \alpha \mathbf{s} \odot \mathbf{g}$	$w_r \leftarrow (1 - \alpha)w_r - \alpha \mathbf{g}$

Table 1. This table compares the steps of our algorithm BayesBiNN to the two existing methods, STE (Bengio et al., 2013) and Bop (Helwegen et al., 2019). Here,  $w_b$  and  $w_r$  denote the binary and real-valued weights. For step 1, where  $w_b$  are obtained from  $w_r$ , STE uses the sign of  $w_r$  while BayesBiNN uses a  $\tanh$  function with a small noise  $\delta$  sampled from a Bernoulli distribution and a temperature parameter  $\tau$ . As Figure 1 (b) shows, as  $\tau$  goes to 0, Step 1 of BayesBiNN becomes equal to that of STE. Step 1 of Bop uses the *hysteresis* function shown in Figure 1 (b) and becomes similar to sign function as the threshold  $\gamma$  goes to 0 (it is flipped but the sign is irrelevant for binary variables). Step 2 is the same for all algorithms. Step 3 of BayesBiNN is very similar to Bop, except that a scaling  $\mathbf{s}$  is used (which is similar to the adaptive learning rate algorithms); see (12) in Section 3.

application. We show an application for continual learning to avoid catastrophic forgetting (Kirkpatrick et al., 2016). To the best of our knowledge, there is no other work on continual learning of BiNNs so far, perhaps because extending existing methods, like STE, for such tasks is not trivial. Overall, our work provides a principled approach for training BiNNs that justifies and extends previous approaches. The code to reproduce the results is available at <https://github.com/team-approx-bayes/BayesBiNN>.

### 1.1. Related Works

There are two main directions on the study of BiNNs: one involves the design of special network architecture tailored to binary operations (Courbariaux et al., 2015; Rastegari et al., 2016; Lin et al., 2017; Bethge et al., 2020) and the other is on the training methods. The latter is the main focus of this paper.

Our algorithm is derived using the Bayesian learning rule recently proposed by Khan & Lin (2017); Khan & Rue (2020), which is obtained by optimizing the Bayesian objective by using natural gradient descent (Amari, 1998; Hoffman et al.,

2013; Khan & Lin, 2017). It is shown in Khan & Rue (2020) that the Bayesian learning rule can be used to derive and justify many existing learning-algorithms in fields such as optimization, Bayesian statistics, machine learning and deep learning. In particular, the Adam optimizer can also be derived as a special case (Khan et al., 2018; Osawa et al., 2019). Our application is yet another example where the rule is used to justify existing algorithms that perform well in practice but whose mechanisms are not well understood.

Instead of using the Bayesian learning rule, it is possible to use other types of variational inference methods, e.g., Shayer et al. (2018); Peters & Welling (2018) used a *variational optimization* approach (Staines & Barber, 2012) along with the local reparameterization trick. The Gumbel-softmax trick (Maddison et al., 2017; Jang et al., 2017) is also used in Louizos et al. (2019) to train BiNNs. However, instead of specifying a Bernoulli distribution over the weights, Louizos et al. (2019) construct a noisy quantizer and their optimization objective is different from ours. Unlike our method, none of these methods (Shayer et al., 2018; Peters & Welling, 2018; Louizos et al., 2019) result in an update similar to either STE or Bop.

## 2. Training Binary Neural Networks (BiNNs)

Given  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ , the goal is to train a neural network  $f_{\mathbf{w}}(\mathbf{x})$  with binary weights  $\mathbf{w} \in \{-1, +1\}^W$ , where  $W$  is the number of weights. The challenge is in optimizing the following discrete optimization objective:

$$\min_{\mathbf{w} \in \{-1, 1\}^W} \sum_{i \in \mathcal{D}} \ell(y_i, f_{\mathbf{w}}(\mathbf{x}_i)), \quad (1)$$

where  $\ell(y_i, \hat{y}_i)$  is a loss function, e.g., cross-entropy loss for the model predictions  $\hat{y}_i := f_{\mathbf{w}}(\mathbf{x}_i)$ . It is clear that binarized weights obtained from pre-trained NNs with real-weights do not minimize (1) and therefore not expected to give good performance. Optimizing the objective with respect to binary weights is difficult since gradient-based methods cannot be directly applied. The gradient of the real-valued weights are not expected to help in the search for the minimum of a discrete objective (Yin et al., 2019).

Despite such theoretical concerns, the Straight-Through Estimator (STE) (Bengio et al., 2013), which utilizes gradient-based methods, works extremely well. There have been many recent works that build upon this method, including BinaryConnect (Courbariaux et al., 2015), Binarized neural networks (Courbariaux et al., 2016), XOR-Net (Rastegari et al., 2016), as well as the most recent MeliusNet (Bethge et al., 2020). The general approach of such methods is shown in Figure 1 in three steps. In step 1, we obtain binary weights  $\mathbf{w}_b$  from the real-valued parameters  $\mathbf{w}_r$ . In step 2, we compute gradients at  $\mathbf{w}_b$  and, in step 3, update  $\mathbf{w}_r$  using the gradients. STE makes a particular choice for the step 1 where a sign function is used to obtain the binary weights from the real-valued weights (see Table 1 for a pseudocode). However, since the gradient of the sign function with respect to  $\mathbf{w}_r$  is zero almost everywhere, it implies that  $\nabla_{\mathbf{w}_r} \approx \nabla_{\mathbf{w}_b}$ . This approximation can be justified in some settings (Yin et al., 2019) but in general the reasons behind its effectiveness are unknown.

Recently Helwegen et al. (2019) proposed a new method that goes against the justification behind STE. They argue that “latent” weights used in STE based methods do not exist. Instead, they provide a new perspective: the sign of each element of  $\mathbf{w}_r$  represents a binary weight while its magnitude encodes some inertia against flipping the sign of the binary weight. With this perspective, they propose the Binary optimizer (Bop) method which keeps track of an exponential moving average (Gardner Jr, 1985) of the gradient  $\mathbf{g}$  during the training process and then decide whether to flip the sign of the binary weights when they exceed a certain threshold  $\gamma$ . The Bop algorithm is shown in Table 1. However, derivation of Bop is also based on intuition and heuristics. It remains unclear why the exponential moving average of the gradient is used in Step 3 and what objective the algorithm is optimizing. The selection of the threshold

$\gamma$  is another difficult choice in the algorithm.

Indeed, Bayesian methods do present a principled way to incorporate the ideas used in both STE and Bop. For example, the idea of “generating” binary weights from real-valued parameters can be thought of as sampling from a discrete distribution with real-valued parameters. In fact, the sign function used in STE is related to the “soft-thresholding” used in machine learning. Despite this there exist no work on Bayesian training of BiNNs that can give an algorithm similar to STE or Bop. In this work, we fix this gap and show that, by using the Bayesian learning rule, we recover a method that justifies some of the steps of STE and Bop, and enable us to extend their application. We will now describe our method.

## 3. BayesBiNN: Binary NNs with Bayes

We will now describe our approach based on a Bayesian formulation of the discrete optimization problem in (1). A Bayesian formulation of a loss-based approach can be written as the following minimization problem with respect to a distribution  $q(\mathbf{w})$  (Zellner, 1988; Bissiri et al., 2016)

$$\mathbb{E}_{q(\mathbf{w})} \left[ \sum_{i=1}^N \ell(y_i, f_{\mathbf{w}}(\mathbf{x}_i)) \right] + \mathbb{D}_{KL}[q(\mathbf{w}) \parallel p(\mathbf{w})], \quad (2)$$

where  $p(\mathbf{w})$  is a prior distribution and  $q(\mathbf{w})$  is the posterior distribution or its approximation. The formulation is general and does not require the loss to correspond to a probabilistic model. When the loss indeed corresponds to a negative log likelihood function, this minimization results in the posterior distribution which is equivalent to Bayes’ rule (Bissiri et al., 2016). When the space of  $q(\mathbf{w})$  is restricted, this results in an approximation to the posterior, which is then equivalent to variational inference (Jordan et al., 1999). For our purpose, this formulation enables us to derive an algorithm that resembles existing methods such as STE and Bop.

### 3.1. BayesBiNN optimizer

To solve the optimization problem (2), the Bayesian learning rule (Khan & Rue, 2020) considers a class of minimal exponential family distributions (Wainwright & Jordan, 2008)

$$q(\mathbf{w}) := h(\mathbf{w}) \exp \left[ \boldsymbol{\lambda}^T \boldsymbol{\phi}(\mathbf{w}) - A(\boldsymbol{\lambda}) \right], \quad (3)$$

where  $\boldsymbol{\lambda}$  is the natural parameter,  $\boldsymbol{\phi}(\mathbf{w})$  is the vector of sufficient statistics,  $A(\boldsymbol{\lambda})$  is the log-partition function, and  $h(\mathbf{w})$  is the base measure. When the prior distribution  $p(\mathbf{w})$  belongs to the same exponential-family as  $q(\mathbf{w})$  in (3), and the base measure  $h(\mathbf{w}) = 1$ , the Bayesian learning uses the following update of the natural parameter (Khan & Lin,

2017; Khan & Rue, 2020)

$$\boldsymbol{\lambda} \leftarrow (1 - \alpha)\boldsymbol{\lambda} - \alpha \left\{ \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q(\boldsymbol{w})} \left[ \sum_{i=1}^N \ell(y_i, f_{\boldsymbol{w}}(\boldsymbol{x}_i)) \right] - \boldsymbol{\lambda}_0 \right\}, \quad (4)$$

where  $\alpha$  is the learning rate,  $\boldsymbol{\mu} = \mathbb{E}_{q(\boldsymbol{w})} [\phi(\boldsymbol{w})]$  is the expectation parameter of  $q(\boldsymbol{w})$ , and  $\boldsymbol{\lambda}_0$  is the natural parameter of the prior distribution  $p(\boldsymbol{w})$ , which is assumed to belong to the same exponential-family as  $q(\boldsymbol{w})$ . The Bayesian learning rule is a natural gradient algorithm (Khan & Lin, 2017; Khan & Rue, 2020). An interesting point of this rule is that the gradient is computed with respect to the expectation parameter  $\boldsymbol{\mu}$ , while the update is performed on the natural parameter  $\boldsymbol{\lambda}$ . This particular choice leads to an update similar to STE for BiNNs, as we show next.

We start by specifying the form of  $p(\boldsymbol{w})$  and  $q(\boldsymbol{w})$ . A priori, we assume that the weights are equally likely to be either  $-1$  or  $+1$ , i.e., the prior  $p(\boldsymbol{w})$  is a (symmetric) Bernoulli distribution with a probability of  $\frac{1}{2}$  for each state. For the posterior approximation  $q(\boldsymbol{w})$ , we use the mean-field (symmetric) Bernoulli distribution

$$q(\boldsymbol{w}) = \prod_{j=1}^W p_j^{\frac{1+w_j}{2}} (1-p_j)^{\frac{1-w_j}{2}}, \quad (5)$$

where  $p_j$  is the probability that  $w_j = +1$ , and  $W$  is the number of parameters. Our goal is to learn the parameters  $p_j$  of the approximations. The Bernoulli distribution defined in (5) is a special case of the *minimal* exponential family distribution, where the corresponding natural and expectation parameters of each weight  $w_i$  are

$$\lambda_j := \frac{1}{2} \log \frac{p_j}{1-p_j}, \quad \mu_j := 2p_j - 1. \quad (6)$$

The natural parameter  $\boldsymbol{\lambda}_0$  of the prior  $p(\boldsymbol{w})$  is therefore  $\mathbf{0}$ . Using these definitions, we can directly apply the Bayesian learning rule to learn the posterior Bernoulli distribution of the binary weights.

In addition to these definitions, we also require the gradient with respect to  $\boldsymbol{\mu}$  to implement the rule (4). A straightforward solution is to use the REINFORCE method (Williams, 1992) which transforms the gradient of the expectation into the expectation of the gradient by using the log-derivative trick, i.e.,

$$\nabla_{\boldsymbol{\mu}} \mathbb{E}_{q(\boldsymbol{w})} [\ell(y, f_{\boldsymbol{w}}(\boldsymbol{x}))] = \mathbb{E}_{q(\boldsymbol{w})} [\ell(y, f_{\boldsymbol{w}}(\boldsymbol{x})) \nabla_{\boldsymbol{\mu}} \log q(\boldsymbol{w})].$$

This method, however, does not use the *minibatch gradient* (the gradient of the loss  $\ell(y, f_{\boldsymbol{w}}(\boldsymbol{x}))$  on a minibatch of examples), which is essential to show the similarity to STE and Bop. The REINFORCE method also suffers from high variance. Due to these reasons, we do not use this method.

Instead, we resort to another reparameterization trick for discrete variables called the Gumbel-softmax trick (Maddison

et al., 2017; Jang et al., 2017), which, as we will see, leads to an update similar to STE/Bop. The idea of Gumbel-softmax trick is to introduce the *concrete distribution* that leads to a *relaxation* of the discrete random variables. Specifically, as shown in Appendix B in (Maddison et al., 2017), for a binary variable  $w_j \in \{0, 1\}$  with  $P(w_j = 1) = p_j$ , we can use the following relaxed variable  $w_r^{\epsilon_j, \tau}(p_j) \in (0, 1)$ :

$$w_r^{\epsilon_j, \tau}(p_j) := \frac{1}{1 + \exp\left(-\frac{2\lambda_j + 2\delta_j}{\tau}\right)}, \quad (7)$$

where  $\tau > 0$  is a temperature parameter,  $\lambda_j := \frac{1}{2} \log \frac{p_j}{1-p_j}$  is the natural parameter, and  $\delta_j$  is defined as follows,

$$\delta_j := \frac{1}{2} \log \frac{\epsilon_j}{1-\epsilon_j}, \quad (8)$$

with  $\epsilon_j \sim \mathcal{U}(0, 1)$  sampled from a uniform distribution. The  $w_r^{\epsilon_j, \tau}(p_j)$  are samples from a Concrete distribution which has a closed-form expression (Maddison et al., 2017):

$$p(w_r^{\epsilon_j, \tau}(p_j)) = \frac{\tau e^{2\lambda} (w_r^{\epsilon_j, \tau}(p_j))^{-\tau-1} (1 - (w_r^{\epsilon_j, \tau}(p_j)))^{-\tau-1}}{\left( e^{2\lambda} (w_r^{\epsilon_j, \tau}(p_j))^{-\tau} + (1 - (w_r^{\epsilon_j, \tau}(p_j)))^{-\tau} \right)^2}. \quad (9)$$

Instead of differentiating the objective with respect to binary variables  $w_j$ , we can differentiate with respect to  $w_r^{\epsilon_j, \tau}(p_j)$ . We will use this to approximate the gradient with respect to  $\boldsymbol{\mu}$  in terms of the minibatch gradient.

In our case, entries  $w_j$  of  $\boldsymbol{w}$  take value in  $\{+1, -1\}$  rather than in  $\{0, 1\}$ . The relaxed version could be obtained by a linear transformation of the concrete variables  $w_r^{\epsilon_j, \tau}(p_j)$  in (7), i.e.,

$$w_b^{\epsilon_j, \tau}(\lambda_j) := 2w_r^{\epsilon_j, \tau}(p_j) - 1 = \tanh((\lambda_j + \delta_j)/\tau), \quad (10)$$

where, unlike (7), we have explicitly written the dependency in terms of  $\lambda_j$  instead of the vector of  $p_j$ . Since  $w_b^{\epsilon_j, \tau}(\lambda_j)$  are continuous, we can differentiate them with respect to  $\mu_j$  by using the chain rule. The lemma below states the result where  $\boldsymbol{w}_b^{\epsilon, \tau}(\boldsymbol{\lambda})$  and  $\boldsymbol{\epsilon}$  denote the vectors of  $w_b^{\epsilon_j, \tau}(\lambda_j)$  and  $\epsilon_j$ , respectively, for all  $j = 1, 2, \dots, W$ .

**Lemma 1** *By using the Gumbel-softmax trick, we get the following approximation in terms of the minibatch gradient:*

$$\nabla_{\boldsymbol{\mu}} \mathbb{E}_{q(\boldsymbol{w})} \left[ \sum_{i=1}^N \ell(y_i, f_{\boldsymbol{w}}(\boldsymbol{x}_i)) \right] \approx \mathbf{s} \odot \mathbf{g}, \quad (11)$$

where

$$\mathbf{g} := \frac{1}{M} \sum_{i \in \mathcal{M}} \nabla_{\boldsymbol{w}_r} \ell(y_i, f_{\boldsymbol{w}_r}(\boldsymbol{x}_i)) \Big|_{\boldsymbol{w}_r = \boldsymbol{w}_b^{\epsilon, \tau}(\boldsymbol{\lambda})}, \quad (12)$$

$$\mathbf{s} := \frac{N(1 - (\boldsymbol{w}_b^{\epsilon, \tau}(\boldsymbol{\lambda}))^2)}{\tau(1 - \tanh(\boldsymbol{\lambda}))^2}, \quad (13)$$

and  $\mathcal{M}$  is a mini-batch of  $M$  examples.

**Proof 1** Using (10), we can first approximate the objective in terms of  $w_b^{\epsilon, \tau}(\lambda)$  and  $\epsilon$ , and then push the gradient with respect to  $\mu$  inside the expectation as shown below:

$$\begin{aligned} \nabla_{\mu} \mathbb{E}_{q(w)} \left[ \sum_{i=1}^N \ell(y_i, f_w(\mathbf{x}_i)) \right] \\ \approx \mathbb{E}_{q(\epsilon)} \left[ \sum_{i=1}^N \nabla_{\mu} \ell(y_i, f_{w_b^{\epsilon, \tau}(\lambda)}(\mathbf{x}_i)) \right]. \end{aligned} \quad (14)$$

Using the chain rule, the  $j$ -th element of the gradient on the right hand side can be obtained as follows:

$$\nabla_{\mu_j} \ell(y_i, f_{w_b^{\epsilon, \tau}(\lambda)}(\mathbf{x}_i)) = \nabla_{w_b^{\epsilon, \tau}} \ell(y_i, f_{w_b^{\epsilon, \tau}(\lambda)}(\mathbf{x}_i)) \frac{dw_b^{\epsilon, \tau}(\lambda_j)}{d\mu_j}. \quad (15)$$

According to the definition of natural parameter and expectation parameter in (6), we have  $\lambda_j = \frac{1}{2} \log \frac{1+\mu_j}{1-\mu_j}$ , therefore after some algebra we can write:

$$\frac{dw_b^{\epsilon, \tau}}{d\mu_j} = \frac{1 - (w_b^{\epsilon, \tau}(\lambda_j))^2}{\tau (1 - \tanh^2(\lambda_j))}. \quad (16)$$

By using a mini-batch  $\mathcal{M}$  of  $M$  examples and one sample  $\epsilon$ , (14)-(16) give us (11).  $\square$

Substituting the result of Lemma 1 into the Bayesian learning rule in (4), we obtain the following update:

$$\lambda \leftarrow (1 - \alpha)\lambda - \alpha [\mathbf{s} \odot \mathbf{g} - \lambda_0]. \quad (17)$$

The resulting optimizer, which we call BayesBiNN, is shown in Table 1, where we assume  $\lambda_0 = \mathbf{0}$  (since the probability of  $w_i = +1$  is  $1/2$  a priori). For the ease of comparison with other methods, the natural parameter  $\lambda$  is replaced with continuous variables  $w_r$ .

At test-time, we can either use the predictions obtained using Monte-Carlo average (which we refer to as the ‘‘mean’’):

$$\hat{p}_k = \frac{1}{C} \sum_{c=1}^C p(y = k | \mathbf{x}, \mathbf{w}^{(c)}), \quad (18)$$

with  $\mathbf{w}^{(c)} \sim q(\mathbf{w})$  and  $C$  is the number of samples, or the predictions obtained by using the mode  $\hat{\mathbf{w}}$  of  $q(\mathbf{w})$ :  $\hat{p}_k = p(y = k | \mathbf{x}, \hat{\mathbf{w}})$ , where  $\hat{\mathbf{w}} = \text{sign}(\tanh(\lambda))$  (which we refer to as the ‘‘mode’’).

### 3.2. Justification of Previous Approaches

In this section, we show how BayesBiNN justifies the steps of STE and Bop. A summary is shown in Table 1. First, BayesBiNN justifies the use of gradient based methods to solve the discrete optimization problem (1). As opposed to

(1), the new objective in (2) is over a continuous parameter  $\lambda$  and thus gradient descent can be used. The underlying principle is similar to stochastic relaxation for non-differentiable optimization (Lemaréchal, 1989; Geman & Geman, 1984), evolution strategies (Huning, 1976), and variational optimization (Staines & Barber, 2012).

Second, some of the algorithmic choices of previous methods such as STE and Bop are justified by BayesBiNN. Specifically, when the temperature  $\tau$  in BayesBiNN is small enough, the  $\tanh(\cdot)$  function in Table 1 behaves like the  $\text{sign}(\cdot)$  function used in STE; see Figure 1 (b). From this perspective, the latent weights  $w_r$  in STE play a similar role as the natural parameter  $\lambda$  of BayesBiNN. In particular, when there is no sampling, i.e.,  $\delta \leftarrow 0$  in BayesBiNN, the two algorithms will become very similar to each other. BayesBiNN justifies the step 1 used in STE by using the Bayesian perspective.

BayesBiNN also justifies step 3 of Bop<sup>1</sup> in Table 1, where an exponential moving average of gradients is used. This is referred to as the momentum term in Bop which plays a similar role as the natural parameter in BayesBiNN. In Helwegen et al. (2019), the momentum is interpreted as a quantity related to *inertia*, which indicates the strength of the state of weights. Since the natural parameter in the binary distribution (5) essentially indicates the strength of the probability being  $-1$  or  $+1$  for each weight, BayesBiNN provides an alternative explanation for Bop.

A recent mirror descent view proposed in Ajanthan et al. (2019b) also interprets the continuous parameters as the dual of the quantized ones. As there is an equivalence between the natural gradient descent and mirror descent (Raskutti & Mukherjee, 2015; Khan & Lin, 2017), the proposed BayesBiNN also provides an interesting perspective on the mirror descent framework for BiNNs training.

### 3.3. Benefits of BayesBiNN

Apart from justifying previous methods, BayesBiNN has several other advantages. First, since its algorithmic form is similar to existing optimizers, it is very easy to implement BayesBiNN by using existing codebases. Second, as a Bayesian method, BayesBiNN provides uncertainty estimates, which can be useful for decision making. The uncertainty obtained using BayesBiNN can enable us to perform continual learning by using the variational continual learning (VCL) framework (Nguyen et al., 2018), as we discuss next.

In continual learning, our goal is to learn the parameters of the neural network from a set of sequentially arriving

<sup>1</sup>Note that the step 3 of Bop in Table 1 is an equivalent but ‘‘flipped’’ version of the one used by Helwegen et al. (2019); see Appendix A for details.

datasets  $\mathcal{D}_t = \{(\mathbf{x}_i, y_i)\}_{i=1}^{N_t}, t = 1, 2, \dots, T$ . While training the  $t$ -th task with dataset  $\mathcal{D}_t$ , we do not have access to the datasets of past tasks, i.e.,  $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_{t-1}$ . Training on task  $\mathcal{D}_t$  using a deep-learning optimizer usually leads to a huge performance loss on the past tasks (Kirkpatrick et al., 2016). The goal of continual learning is to fix such catastrophic forgetting of the past.

For full-precision networks, a common approach to solve this problem is to use weight-regularization, e.g., the *elastic weight consolidation* (EWC) method (Kirkpatrick et al., 2016) uses a Fisher information matrix to regularize the weights:

$$\sum_{i \in \mathcal{D}_t} \ell(y_i, f_{\mathbf{w}}(\mathbf{x}_i)) + \varepsilon (\mathbf{w} - \mathbf{w}_{t-1})^T \mathbf{F}_t (\mathbf{w} - \mathbf{w}_{t-1}), \quad (19)$$

where  $\varepsilon$  is a regularization parameter and  $\mathbf{F}_t$  is the Fisher information matrix at  $\mathbf{w}_t$ . The hope is to keep the new weights close to the old weights, but in a discrete optimization problem, it is impossible to characterize such closeness using a quadratic regularizer as above. Therefore, it is unclear why such a regularizer will be useful. In addition, the use of the Fisher information matrix  $\mathbf{F}_t$  typically assumes that the weights are continuous and the matrix does not provide a meaningful quantity for discrete weights. For these reasons, extending existing approaches such as STE and Bop to continual learning is a nontrivial task.

Fortunately, for BayesBiNN, this is very easy because the objective function is well-defined. We use the VCL framework (Nguyen et al., 2018) where we regularize the distributions instead of the weights. The Kullback-Leibler term is used as the regularizer. Denoting by  $q_{t-1}(\mathbf{w})$  the posterior distribution at task  $t - 1$ , we can replace the prior distribution  $p(\mathbf{w})$  in (2) by  $q_{t-1}(\mathbf{w})$ :

$$\mathbb{E}_{q_t(\mathbf{w})} \left[ \sum_{i \in \mathcal{D}_t} \ell(y_i, f_{\mathbf{w}}(\mathbf{x}_i)) \right] + \mathbb{D}_{KL}[q_t(\mathbf{w}) \| q_{t-1}(\mathbf{w})]. \quad (20)$$

This leads to a slight modification of the update in BayesBiNN where the prior natural parameter  $\lambda_0$  in (17) is replaced by the natural parameter  $\lambda_{t-1}$  of  $q_{t-1}(\mathbf{w})$ . The new update of the natural parameter  $\lambda_t$  of  $q_t(\mathbf{w})$  is shown below:

$$\lambda_t \leftarrow (1 - \alpha)\lambda_t - \alpha [\mathbf{s} \odot \mathbf{g} - \lambda_{t-1}]. \quad (21)$$

By using a posterior approximation  $q(\mathbf{w})$  and a well-defined objective, BayesBiNN enables the application of STE/Bop like methods to such challenging continual learning problems.

## 4. Experimental Results

In this section, we present numerical experiments to demonstrate the performance of BayesBiNN on both synthetic and

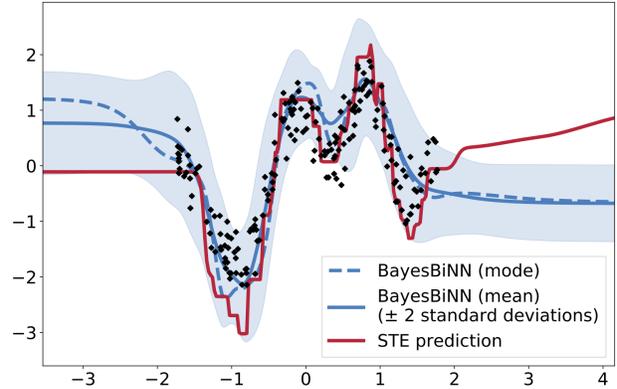


Figure 2. Regression on Snelson dataset (Snelson & Ghahramani, 2005). BayesBiNN (mean) gives much smoother curve than the STE. Uncertainty is low in areas with plenty of data points.

real image data for different kinds of neural network architectures. We also show an application of BayesBiNN to continual learning. The code to reproduce the results is available at <https://github.com/team-approx-bayes/BayesBiNN>.

### 4.1. Synthetic Data

First, we present visualizations on toy regression and binary classification problems. The STE (Bengio et al., 2013) algorithm is used as a baseline for which we employ Adam for training. We use a multi-layer perceptron (MLP) with two hidden layers of 64 units and tanh activation functions.

**Regression** In Figure 2, we show results for regression on the Snelson dataset (Snelson & Ghahramani, 2005). For this experiment, we add a Batch normalization (BN) (Ioffe & Szegedy, 2015) layer (but no learned gain or bias terms) after the last fully connected layer. As seen in Figure 2, predictions obtained using ‘BayesBiNN (mean)’ gives much smoother curves than STE, as expected. Uncertainty is lower in the areas with little noise and plenty of data points compared to areas with no data. Experimental details of the training process are provided in the Appendix B.1.

**Classification** Figure 3 shows STE and BayesBiNN on the two moons dataset (Moons) with 100 data points in each class. STE (the leftmost figure) gives a point estimate of the weights and results in a fairly deterministic classifier. When using the mode of the BayesBiNN distribution (the middle figure), the results are similar, with the fit being slightly worse but overall less overconfident, especially in the regions with no data. Using the mean over 10 samples drawn from the posterior distribution  $q(\mathbf{w})$  (the rightmost figure), we get much better uncertainty estimates as we move away from the data. Experimental details of the training process

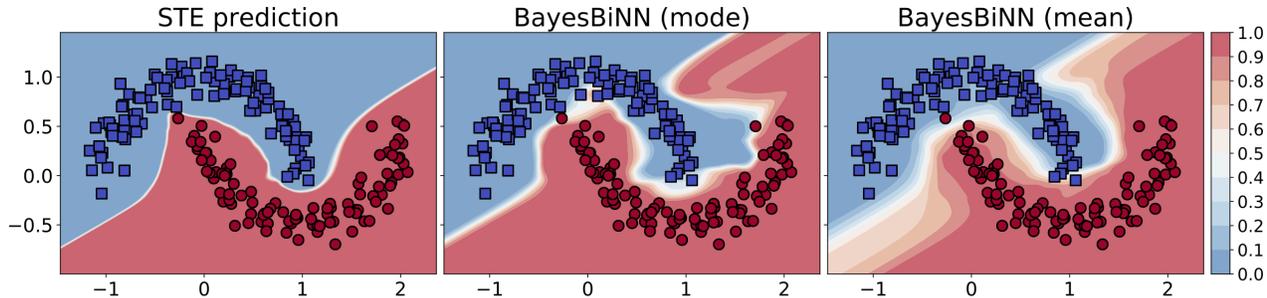


Figure 3. Classification on the two moons dataset with different optimizers. From left to right: STE, BayesBiNN using the mode, BayesBiNN using the predictive mean from 10 posterior Monte Carlo samples. STE is more overconfident than BayesBiNN, and BayesBiNN (mean) gives reasonable uncertainty in regions further away from the data.

Table 2. Test accuracy of different optimizers for MNIST, CIFAR-10 and CIFAR-100 (Averaged over 5 runs). In all the three benchmark datasets, BayesBiNN achieves similar performance as STE Adam and closely approaches the performance of full-precision networks.

Optimizer	MNIST	CIFAR-10	CIFAR-100
STE Adam	<b>98.85</b> $\pm$ 0.09 %	<b>93.55</b> $\pm$ 0.15 %	72.89 $\pm$ 0.21 %
Bop	98.47 $\pm$ 0.02 %	93.00 $\pm$ 0.11 %	69.58 $\pm$ 0.15 %
PMF	98.80 $\pm$ 0.06 %	91.43 $\pm$ 0.14 %	70.45 $\pm$ 0.25 %
<b>BayesBiNN (mode)</b>	<b>98.86</b> $\pm$ 0.05 %	<b>93.72</b> $\pm$ 0.16 %	<b>73.68</b> $\pm$ 0.31 %
<b>BayesBiNN (mean)</b>	<b>98.86</b> $\pm$ 0.05 %	<b>93.72</b> $\pm$ 0.15 %	<b>73.65</b> $\pm$ 0.41 %
Full-precision	99.01 $\pm$ 0.06 %	93.90 $\pm$ 0.17 %	74.83 $\pm$ 0.26 %

are provided in the Appendix B.1.

## 4.2. Image Classification on Real Datasets

We now present results on three benchmark real datasets widely used for image classification: MNIST (LeCun & Cortes, 2010), CIFAR-10 (Krizhevsky & Hinton, 2009) and CIFAR-100 (Krizhevsky & Hinton, 2009). We compare to three other optimizers<sup>2</sup>: STE (Bengio et al., 2013) using Adam with weight clipping and gradient clipping (Courbariaux et al., 2015; Maddison et al., 2017; Alizadeh et al., 2019); latent-free Bop (Helweggen et al., 2019); and the proximal mean-field (PMF) (Ajanthan et al., 2019a). An additional comparison with the LR-net method of Shayer et al. (2018) is given in Appendix B.3. For a fair comparison, we keep all conditions the same except for the optimization methods themselves. For our proposed BayesBiNN, we report results using both the mode and the mean. For all the experiments, standard categorical cross-entropy loss is used and we take 10% of the training set for validation and report the best accuracy on the test set corresponding to the highest validation accuracy achieved during training.

For MNIST, we use a multilayer perceptron (MLP) with

<sup>2</sup>We use Bop code available at <https://github.com/plumerai/rethinking-bnn-optimization> and PMF code available at <https://github.com/tajanthan/pmf>.

three hidden layers with 2048 units and rectified linear units (ReLU) (Alizadeh et al., 2019) activations. Both Batch normalization (BN)<sup>3</sup> (Ioffe & Szegedy, 2015) and dropout are used. No data augmentation is performed. For CIFAR-10 and CIFAR-100, we use the BinaryConnect CNN network in Alizadeh et al. (2019), which is a VGG-like structure similar to the one used in Helweggen et al. (2019). Standard data augmentation is used (Graham, 2014), where 4 pixels are padded on each side, a random  $32 \times 32$  crop is applied, followed by a random horizontal flip. Note that no ZCA whitening is used as in Courbariaux et al. (2015); Alizadeh et al. (2019). The details of the experimental setting, including the detailed network architecture and values of all hyper-parameters, are provided in Appendix B.2 in the supplementary material.

As shown in Table 2, the proposed BayesBiNN achieves similar performances (slightly better for CIFAR-100) as STE Adam in all the three benchmark datasets and approaches the performance of full-precision DNNs. The detailed results, such as the train/validation accuracy as well as the training curves are provided in Appendix B.2 in the supplementary material.

<sup>3</sup>Here the parameters of BN layers are not learned. However, they could also be learned by applying a conventional optimizer such as Adam separately, which is easy to implement.

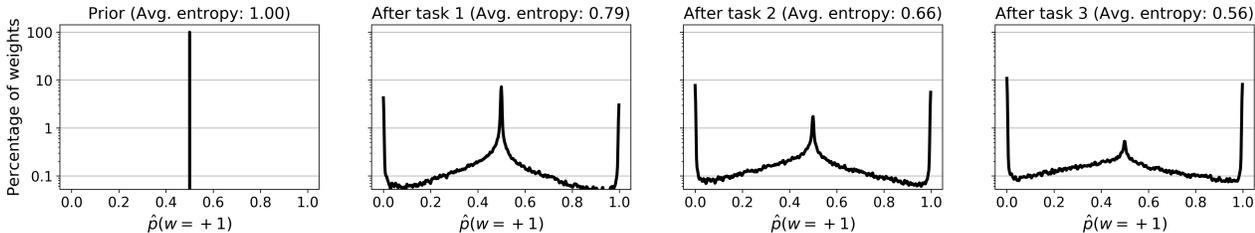


Figure 4. Evolution of the distribution over the binary weights during the learning process in continual learning. The histogram of the weight probabilities  $\hat{p}(w_j = +1)$  for all  $j$  is shown after learning on different tasks. At the very beginning, all the weights are equal to -1 or +1 with prior probability 0.5 and thus have maximum average entropy 1.0. As the number of learned tasks increases, the distribution spreads and becomes flatter, implying that the average entropy of the binary weights decreases, i.e., the weights of BiNNs become more and more deterministic.

### 4.3. Continual learning with binary neural networks

We now show an application of BayesBiNN to continual learning. We consider the popular benchmark of permuted MNIST (Goodfellow et al., 2013; Kirkpatrick et al., 2016; Nguyen et al., 2018; Zenke et al., 2017), where each dataset  $\mathcal{D}_t$  consists of labeled MNIST images whose pixels are permuted randomly. Similar to Nguyen et al. (2018), we use a fully connected single-head network with two hidden layers containing 100 hidden units with ReLU activations. No coresets are used. The details of the experiment, e.g., the network architecture and values of hyper-parameters, are provided in Appendix B.4.

As shown in Figure 5, using BayesBiNN with the posterior approximation  $q_{t-1}(\mathbf{w})$  as the prior for task  $t$  (red solid line), we achieve significant improvements in overcoming catastrophic forgetting of the past. When the prior is fixed to be  $\lambda_0 = \mathbf{0}$  (blue dotted line), the network performs badly on the past tasks, e.g., in the top row, after trained on task 2 and task 3, the network performs badly on the previous task 1. The reason for better performance when using the posterior approximation  $q_{t-1}(\mathbf{w})$  as the prior for task  $t$  is directly related to the uncertainty estimated by the posterior approximation  $q_{t-1}(\mathbf{w})$ . To visualize the uncertainty, Figure 4 shows the histogram of the weight probabilities  $\hat{p}(w_j = +1)$  for all  $j$ . The prior probability, shown in the first plot, is set to  $\frac{1}{2}$  for all weights (entropy is 1.0). As we train on more tasks, the uncertainty decreases and the weights of BiNNs become more and more deterministic (the distribution spreads and becomes flatter). As desired, with more data, the network reduces the entropy of the distribution and the uncertainty is useful to perform continual learning.

## 5. Conclusion

Binary neural networks (BiNNs) are computation-efficient and hardware-friendly, but their training is challenging since in theory it involves a discrete optimization problem. How-

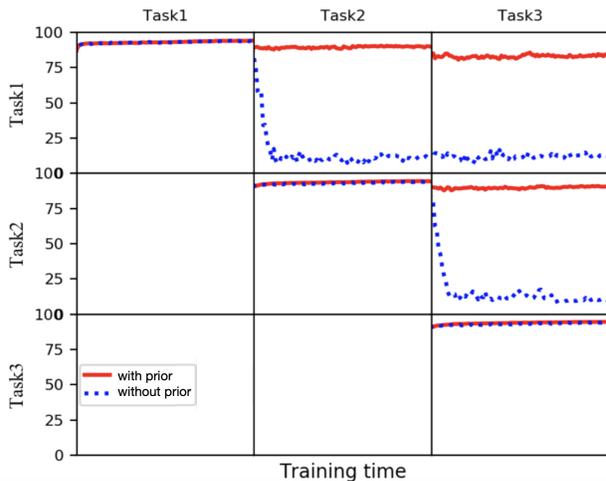


Figure 5. Test accuracy curves of continual learning for BiNNs using BayesBiNN on permuted MNIST. The neural network is trained with or without prior for BayesBiNN, respectively. Specifically, with prior (red solid line) indicates using BayesBiNN with the posterior approximation  $q_{t-1}(\mathbf{w})$  as the prior for task  $t$  while without prior (blue dotted line) indicates using BayesBiNN with fixed prior where  $\lambda_0 = \mathbf{0}$ . The test accuracy on the test set is averaged over 5 random runs. The X-axis shows the training time (epochs) and Y-axis shows the average test accuracy of different tasks as the training time increases.

ever, some gradient-based methods such as the STE work quite well in practice despite ignoring the discrete nature of the optimization problem, which is surprising and there is a lack of principled justification of their success. In this paper, we proposed a principled approach to train the binary neural networks using the Bayesian learning rule. The resulting optimizer, which we call BayesBiNN, not only justifies some of the algorithmic choices made by existing methods such as the STE and Bop but also facilitates the extensions of them, e.g., enabling uncertainty estimation for continual learning to avoid the catastrophic forgetting problem for binary neural networks.

## Acknowledgements

We would like to thank the anonymous reviewers for their valuable comments. X. Meng would like to thank Milad Alizadeh (Oxford University) and Koen Helwegen (Plumerai Research) for useful discussions. We are thankful for the RAIDEN computing system and its support team at the RIKEN Center for Advanced Intelligence Project which we used extensively for our experiments.

## References

- Ajanthan, T., Dokania, P. K., Hartley, R., and Torr, P. H. S. Proximal mean-field for neural network quantization. In *IEEE International Conference on Computer Vision*, pp. 4871–4880, 2019a.
- Ajanthan, T., Gupta, K., Torr, P. H., Hartley, R., and Dokania, P. K. Mirror descent view for neural network quantization. *arXiv preprint arXiv:1910.08237*, 2019b.
- Alizadeh, M., Fernández-Marqués, J., Lane, N. D., and Gal, Y. An empirical study of binary neural networks’ optimisation. *ICLR*, 2019.
- Amari, S.-I. Natural gradient works efficiently in learning. *Neural computation*, 10(2):251–276, 1998.
- Bengio, Y., Léonard, N., and Courville, A. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.
- Bethge, J., Bartz, C., Yang, H., Chen, Y., and Meinel, C. Meliusnet: Can binary neural networks achieve mobilenet-level accuracy? *arXiv preprint arXiv:2001.05936*, 2020.
- Bissiri, P. G., Holmes, C. C., and Walker, S. G. A general framework for updating belief distributions. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 78(5):1103–1130, 2016.
- Courbariaux, M., Bengio, Y., and David, J.-P. Binaryconnect: Training deep neural networks with binary weights during propagations. In *Advances in neural information processing systems*, pp. 3123–3131, 2015.
- Courbariaux, M., Hubara, I., Soudry, D., El-Yaniv, R., and Bengio, Y. Binarized neural networks: Training deep neural networks with weights and activations constrained to+ 1 or-1. *arXiv preprint arXiv:1602.02830*, 2016.
- Gardner Jr, E. S. Exponential smoothing: The state of the art. *Journal of forecasting*, 4(1):1–28, 1985.
- Geman, S. and Geman, D. Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. *IEEE Transactions on pattern analysis and machine intelligence*, (6):721–741, 1984.
- Goodfellow, I. J., Mirza, M., Xiao, D., Courville, A., and Bengio, Y. An empirical investigation of catastrophic forgetting in gradient-based neural networks. *ICLR*, 2013.
- Graham, B. Spatially-sparse convolutional neural networks. *arXiv preprint arXiv:1409.6070*, 2014.
- Helwegen, K., Widdicombe, J., Geiger, L., Liu, Z., Cheng, K.-T., and Nusselder, R. Latent weights do not exist: Rethinking binarized neural network optimization. *NeurIPS*, 2019.
- Hoffman, M. D., Blei, D. M., Wang, C., and Paisley, J. Stochastic variational inference. *The Journal of Machine Learning Research*, 14(1):1303–1347, 2013.
- Huning, A. Evolutionsstrategie. optimierung technischer systeme nach prinzipien der biologischen evolution, 1976.
- Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *ICML*, 2015.
- Jang, E., Gu, S., and Poole, B. Categorical reparameterization with Gumbel-softmax. *ICLR*, 2017.
- Jordan, M. I., Ghahramani, Z., Jaakkola, T. S., and Saul, L. K. An introduction to variational methods for graphical models. *Machine learning*, 37(2):183–233, 1999.
- Khan, M. E. and Lin, W. Conjugate-computation variational inference: Converting variational inference in non-conjugate models to inferences in conjugate models. *AIS-TATS*, 2017.
- Khan, M. E. and Rue, H. Learning-Algorithms from Bayesian principles. 2020. [https://emtiyaz.github.io/papers/learning\\_from\\_bayes.pdf](https://emtiyaz.github.io/papers/learning_from_bayes.pdf).
- Khan, M. E., Nielsen, D., Tangkaratt, V., Lin, W., Gal, Y., and Srivastava, A. Fast and scalable Bayesian deep learning by weight-perturbation in Adam. *ICML*, 2018.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *ICLR*, 2015.
- Kirkpatrick, J. N., Pascanu, R., Rabinowitz, N. C., Veness, J., Desjardins, G., Rusu, A. A., Milan, K., Quan, J., Ramalho, T., Grabska-Barwinska, A., Hassabis, D., Clopath, C., Kumaran, D., and Hadsell, R. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences of the United States of America*, 114 13:3521–3526, 2016.

- Krizhevsky, A. and Hinton, G. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- Langley, P. Crafting papers on machine learning. In Langley, P. (ed.), *Proceedings of the 17th International Conference on Machine Learning (ICML 2000)*, pp. 1207–1216, Stanford, CA, 2000. Morgan Kaufmann.
- LeCun, Y. and Cortes, C. MNIST handwritten digit database. 2010. URL <http://yann.lecun.com/exdb/mnist/>.
- Lemaréchal, C. Nondifferentiable optimization. *Handbooks in operations research and management science*, 1:529–572, 1989.
- Lin, X., Zhao, C., and Pan, W. Towards accurate binary convolutional neural network. In *Advances in Neural Information Processing Systems*, pp. 345–353, 2017.
- Louizos, C., Reisser, M., Blankevoort, T., Gavves, E., and Welling, M. Relaxed quantization for discretized neural networks. *ICLR*, 2019.
- Maddison, C. J., Mnih, A., and Teh, Y. W. The concrete distribution: A continuous relaxation of discrete random variables. *ICLR*, 2017.
- Mishra, A., Nurvitadhi, E., Cook, J. J., and Marr, D. WRPN: wide reduced-precision networks. *ICLR*, 2018.
- Moons, T. Two moons datasets description. [https://scikit-learn.org/stable/modules/generated/sklearn.datasets.make\\_moons.html](https://scikit-learn.org/stable/modules/generated/sklearn.datasets.make_moons.html).
- Nguyen, C. V., Li, Y., Bui, T. D., and Turner, R. E. Variational continual learning. *ICLR*, 2018.
- Osawa, K., Swaroop, S., Khan, M. E. E., Jain, A., Eschenhagen, R., Turner, R. E., and Yokota, R. Practical deep learning with Bayesian principles. In *Advances in Neural Information Processing Systems*, pp. 4289–4301, 2019.
- Peters, J. W. and Welling, M. Probabilistic binary neural networks. *arXiv preprint arXiv:1809.03368*, 2018.
- Raskutti, G. and Mukherjee, S. The information geometry of mirror descent. *IEEE Transactions on Information Theory*, 61(3):1451–1457, 2015.
- Rastegari, M., Ordonez, V., Redmon, J., and Farhadi, A. Xnor-net: Imagenet classification using binary convolutional neural networks. In *European Conference on Computer Vision*, pp. 525–542. Springer, 2016.
- Shayer, O., Levi, D., and Fetaya, E. Learning discrete weights using the local reparameterization trick. *ICLR*, 2018.
- Snelson, E. and Ghahramani, Z. Sparse Gaussian processes using pseudo-inputs. In *Proceedings of the 18th International Conference on Neural Information Processing Systems, NIPS05*, pp. 12571264, Cambridge, MA, USA, 2005. MIT Press.
- Staines, J. and Barber, D. Variational optimization. *arXiv preprint arXiv:1212.4507v2*, 2012.
- Wainwright, M. J. and Jordan, M. I. Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, 1–2:1–305, 2008.
- Williams, R. J. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.
- Yin, P., Lyu, J., Zhang, S., Osher, S., Qi, Y., and Xin, J. Understanding straight-through estimator in training activation quantized neural nets. *ICLR*, 2019.
- Zellner, A. Optimal information processing and bayes’s theorem. *The American Statistician*, 42(4):278–280, 1988.
- Zenke, F., Poole, B., and Ganguli, S. Continual learning through synaptic intelligence. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 3987–3995. JMLR. org, 2017.

---

# Training Binary Neural Networks using the Bayesian Learning Rule: Appendix

---

## A. Two equivalent forms of hysteresis function in Bop

The original update rule and the corresponding definition of the *hysteresis* function  $\text{hyst}(\cdot)$  in Bop are (Helweggen et al., 2019)

$$\mathbf{w}_r \leftarrow (1 - \alpha)\mathbf{w}_r + \alpha\mathbf{g}, \quad (22)$$

$$\begin{aligned} y &= \text{hyst1}(w_r, w_b, \gamma) \\ &\equiv \begin{cases} -w_b & \text{if } |w_r| > \gamma \text{ \& } \text{sign}(w_r) = \text{sign}(w_b), \\ w_b & \text{otherwise.} \end{cases} \end{aligned} \quad (23)$$

One could obtain an alternative update rule  $\mathbf{w}_r \leftarrow (1 - \alpha)\mathbf{w}_r - \alpha\mathbf{g}$ , as shown in Step 3 of Bop in Table 1. In this case, the update rule and the corresponding *hysteresis* function are as follows

$$\mathbf{w}_r \leftarrow (1 - \alpha)\mathbf{w}_r - \alpha\mathbf{g}, \quad (24)$$

$$\begin{aligned} y &= \text{hyst2}(w_r, w_b, \gamma) \\ &\equiv \begin{cases} -w_b & \text{if } |w_r| > \gamma \text{ \& } \text{sign}(w_r) = -\text{sign}(w_b), \\ w_b & \text{otherwise.} \end{cases} \end{aligned} \quad (25)$$

It could be easily verified that the above two update rules with two different representations of the *hysteresis* function are equivalent to each other: The only difference between (22) and (24) is the sign before the gradient  $\mathbf{g}$ , i.e., the  $\mathbf{w}_r$  in (22) is an exponential moving average (Gardner Jr, 1985) of  $\mathbf{g}$  while in (24) it is an exponential moving average of  $-\mathbf{g}$ . Such difference is compensated by the difference between (23) and (25). The corresponding curve of  $y = \text{hyst1}(w_r, w_b, \gamma)$  is simply a upside-down flipped version of  $y = \text{hyst2}(w_r, w_b, \gamma)$ , which is shown in the rightmost figure in Figure 1 (b).

## B. Experimental details

In this section we list the details for all experiments shown in the main text.

Note that after training BiNNs with BayesBiNN, there are two ways to perform inference during test time:

(1). **Mean:** One method is to use the predictive mean, where we use Monte Carlo sampling to compute the predictive probabilities for each test sample  $\mathbf{x}_j$  as follows

$$\hat{p}_{j,k} \approx \frac{1}{C} \sum_{c=1}^C p(y_j = k | \mathbf{x}_j, \mathbf{w}^{(c)}), \quad (26)$$

where  $\mathbf{w}^{(c)} \sim q(\mathbf{w})$  are samples from the Bernoulli distributions with the natural parameters  $\boldsymbol{\lambda}$  obtained by BayesBiNN.

(2). **Mode:** The other way is simply to use the mode of the posterior distribution  $q(\mathbf{w})$ , i.e., the sign value of the posterior mean, i.e.,  $\hat{\mathbf{w}} = \text{sign}(\tanh(\boldsymbol{\lambda}))$ , to make predictions, which will be denoted as  $C = 0$ .

### B.1. Synthetic Data

**Binary Classification** We used the Two Moons dataset with 100 data points in each class and added Gaussian noise with standard deviation 0.1 to each point. We trained a Multilayer Perceptron (MLP) with two hidden layers of 64 units and tanh

activation functions for 3000 epochs, using Cross Entropy as the loss function. Additional train and test settings with respect to the optimizers are detailed in Table 3. The learning rate  $\alpha$  was decayed at fixed epochs by the specified learning rate decay rate. For the STE baseline, we used the Adam optimizer with standard settings.

Table 3. Train settings for the binary classification experiment using the Two Moons dataset.

Setting	BayesBiNN	STE
Learning rate $\alpha$	$10^{-3}$	$10^{-1}$
Learning rate decay	0.1	0.1
Learning rate decay epochs	[1500, 2500]	[1500, 2500]
Momentum(s) $\beta$	0.99	0.9, 0.999
MC train samples $S$	5	-
MC test samples $C$	0/10	-
Temperature $\tau$	1	-
Prior $\lambda_0$	<b>0</b>	-
Initialization $\lambda$	$\pm 15$ randomly	-

**Regression** We used the Snelson dataset (Snelson & Ghahramani, 2005) with 200 data points to train a regression model. Similar to the Binary Classification experiment, we used an MLP with two hidden layers of 64 units and tanh activation functions, but trained it for 5000 epochs using Mean Squared Error as the loss function. Additionally, we added a batch normalization layer (without learned gain or bias terms) after the last fully connected layer. The learning rate is adjusted after every epoch to slowly anneal from an initial learning rate  $\alpha_0$  to a target learning rate  $\alpha_T$  at the maximum epoch  $T$  using

$$\alpha_{t+1} = \alpha_t \left( \frac{\alpha_T}{\alpha_0} \right)^{-T}. \tag{27}$$

The learning rates and other train and test settings are detailed in Table 4.

Table 4. Train settings for the regression experiment using the Snelson dataset (Snelson & Ghahramani, 2005).

Setting	BayesBiNN	STE
Learning rate start $\alpha_0$	$10^{-4}$	$10^{-1}$
Learning rate end $\alpha_T$	$10^{-5}$	$10^{-1}$
Momentum(s) $\beta$	0.99	0.9, 0.999
MC train samples $S$	1	-
MC test samples $C$	0/10	-
Temperature $\tau$	1	-
Prior $\lambda_0$	<b>0</b>	-
Initialization $\lambda$	$\pm 10$ randomly	-

## B.2. MNIST, CIFAR-10 and CIFAR-100

In this section, three well-known image datasets are considered, namely the MNIST, CIFAR-10 and CIFAR-100 datasets. We compare the proposed BayesBiNN with four other popular algorithms, STE Adam, Bop and PMF for BiNNs as well as standard Adam for full-precision weights. For dataset and algorithm specific settings, see Table 9.

**MNIST** All algorithms have been trained using the same MLP detailed in Table 5 on mini-batches of size 100, for a maximum of 500 epochs. The loss used was Categorical Cross Entropy. We split the original training data into 90% train and 10% validation data and no data augmentation except normalization has been done. We report the best accuracy (averaged over 5 random runs) on the test set corresponding to the highest validation accuracy achieved during training (we do not retrain using the validation set). Note that we tune the hyper-parameters such as learning rate for all the methods including the baselines. The search space for the learning rate is set to be  $[10^{-2}, 3 \cdot 10^{-3}, 10^{-3}, 3 \cdot 10^{-4}, 10^{-4}, 3 \cdot 10^{-5}, 10^{-5}, 10^{-6}]$

for all methods. Moreover, Table 6 and Table 7 shows the results of MNIST with BayesBiNN for different choices of learning rate and temperature.

Table 5. The MLP architecture used in all MNIST experiments, adapted from (Alizadeh et al., 2019).

Dropout (p = 0.2)
Fully Connected Layer (units = 2048, bias = False)
ReLU
Batch Normalization Layer (gain = 1, bias = 0)
Dropout (p = 0.2)
Fully Connected Layer (units = 2048, bias = False)
ReLU
Batch Normalization Layer (gain = 1, bias = 0)
Dropout (p = 0.2)
Fully Connected Layer (units = 2048, bias = False)
ReLU
Batch Normalization Layer (gain = 1, bias = 0)
Dropout (p = 0.2)
Fully Connected Layer (units = 2048, bias = False)
Batch Normalization Layer (gain = 1, bias = 0)
Softmax

Table 6. Test accuracy of MNIST for different initial learning rates. The temperature is  $10^{-10}$ . Results are averaged over 5 random runs.

Learning rate	$10^{-1}$	$3 \cdot 10^{-3}$	$10^{-3}$	$3 \cdot 10^{-4}$
Training Accuracy	$99.46 \pm 0.15 \%$	$99.58 \pm 0.16 \%$	$99.67 \pm 0.09 \%$	$99.76 \pm 0.09 \%$
Validation Accuracy	$98.90 \pm 0.14 \%$	$98.94 \pm 0.17 \%$	$98.96 \pm 0.13 \%$	$98.97 \pm 0.12 \%$
Test Accuracy	$98.73 \pm 0.11 \%$	$98.81 \pm 0.07 \%$	$98.83 \pm 0.05 \%$	$98.84 \pm 0.08 \%$
Learning rate	$10^{-4}$	$3 \cdot 10^{-5}$	$10^{-5}$	$10^{-6}$
Training Accuracy	$99.85 \pm 0.05 \%$	$99.83 \pm 0.06 \%$	$99.76 \pm 0.09 \%$	$99.78 \pm 0.03 \%$
Validation Accuracy	$99.02 \pm 0.13 \%$	$99.02 \pm 0.13 \%$	$99.04 \pm 0.11 \%$	$99.02 \pm 0.17 \%$
Test Accuracy	$98.86 \pm 0.05 \%$	$98.86 \pm 0.05 \%$	$98.84 \pm 0.08 \%$	$98.85 \pm 0.05 \%$

**CIFAR-10 and CIFAR-100** We trained all algorithms on the Convolutional Neural Network (CNN) architecture detailed in Table 8 on mini-batches of size 50, for a maximum of 500 epochs. The loss used was Categorical Cross Entropy. We split the original training data into 90% train and 10% validation data. For data augmentation during training, the images were normalized, a random  $32 \times 32$  crop was selected from a  $40 \times 40$  padded image and finally a random horizontal flip was applied. In the same manner as Osawa et al. (2019), we consider such data augmentation as effectively increasing the dataset size by a factor of 10 (4 images for each corner, and one central image, and the horizontal flipping step further doubles the dataset size, which gives a total factor of 10). We report the best accuracy (averaged over 5 random runs) on the test set corresponding to the highest validation accuracy achieved during training. In addition, we tune the hyper-parameters, such as the learning rate, for all the methods including the baselines. The search space for the learning rate is set to be  $[10^{-2}, 3 \cdot 10^{-3}, 10^{-3}, 3 \cdot 10^{-4}, 10^{-4}, 3 \cdot 10^{-5}, 10^{-5}, 10^{-6}]$  for all methods.

### B.3. Comparison with LR-net

We also compare the proposed BayesBiNN with the LR-net method in Shayer et al. (2018) for MNIST and CIFAR-10. As the code for the LR-net is not open-source, we performed experiments with BayesBiNN following the same experimental

Table 7. Test accuracy of MNIST for different temperatures. The initial learning rate is  $10^{-4}$ . Results are averaged over 5 random runs.

Temperature	$10^{-3}$	$10^{-4}$	$10^{-5}$	$10^{-6}$	$10^{-7}$
Training Accuracy	$89.25 \pm 0.22$ %	$87.55 \pm 0.50$ %	$90.22 \pm 0.42$ %	$97.37 \pm 0.13$ %	$98.27 \pm 0.10$ %
Validation Accuracy	$90.06 \pm 1.04$ %	$90.28 \pm 0.43$ %	$93.35 \pm 0.48$ %	$98.10 \pm 0.17$ %	$98.55 \pm 0.16$ %
Test Accuracy	$90.40 \pm 0.97$ %	$90.72 \pm 0.42$ %	$93.67 \pm 0.50$ %	$98.01 \pm 0.05$ %	$98.41 \pm 0.10$ %
Learning rate	$10^{-8}$	$10^{-9}$	$10^{-10}$	$10^{-11}$	$10^{-12}$
Training Accuracy	$99.48 \pm 0.08$ %	$99.75 \pm 0.14$ %	$99.85 \pm 0.05$ %	$99.81 \pm 0.04$ %	$99.82 \pm 0.07$ %
Validation Accuracy	$98.92 \pm 0.13$ %	$99.00 \pm 0.13$ %	$99.02 \pm 0.14$ %	$99.02 \pm 0.12$ %	$99.02 \pm 0.13$ %
Test Accuracy	$98.82 \pm 0.05$ %	$98.81 \pm 0.08$ %	$98.86 \pm 0.05$ %	$98.86 \pm 0.06$ %	$98.84 \pm 0.04$ %

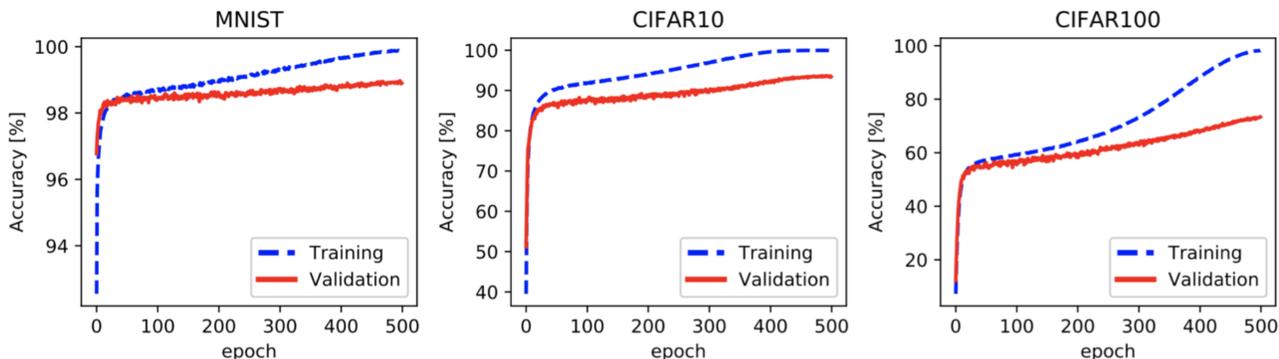


Figure 6. Training/Validation accuracy for MNIST, CIFAR-10 and CIFAR100 with BayesBiNN optimizer (Averaged over 5 runs).

settings in Shayer et al. (2018) and then compared the results with the reported results in their paper. In specific, the network architectures for MNIST and CIFAR-10 are the same as Shayer et al. (2018), except that we added BN after the FC layers. However, we kept all layers binary and did not learn the BN parameters, nor did we use dropout as in Shayer et al. (2018). The dataset pre-processing follows the same settings in Shayer et al. (2018) and is similar to that described in subsection 4.2, except that there is no split of the training set into training and validation sets. As a result, as in Shayer et al. (2018), we report the test accuracies after 190 epochs and 290 epochs for MNIST and CIFAR-10, respectively. Note that the hyper-parameter settings of BayesBiNN are the same as those in Table 9 for MNIST and CIFAR-10. The results are shown in Table 11. The proposed BayesBiNN achieves similar performance (slightly better for CIFAR-10) to the LR-net. Note that the LR-net method used pre-trained models to initialize the weights of BiNNs, while BayesBiNN trained BiNNs from scratch without using pre-trained models.

#### B.4. Continual learning with binary neural networks

For the continual learning experiment, we used a three-layer MLP, detailed in Table 12, and trained it using the Categorical Cross Entropy loss. Specific training parameters are given in Table 13. There is no split of the original MNIST training data in the continual learning case. No data augmentation except normalization has been performed.

### C. Author Contributions Statement

M.E.K. conceived the idea of training Binary neural networks using the Bayesian learning rule. X.M. derived the BayesBiNN algorithm, studied its connections to STE and Bop, and wrote the first proof-of-concept experiments. R.B. fixed a few issue with the original implementation and re-organized the PyTorch code. R.B. also designed and performed the experiments on synthetic data presented in Section 4.1. X.M. did most of the experiments with some help from R.B. All the authors were involved in writing, revising and proof-reading the paper.

Table 8. The CNN architecture used in all CIFAR-10 and CIFAR-100 experiments, inspired by VGG and used in Alizadeh et al. (2019).

Convolutional Layer (channels = 128, kernel-size = $3 \times 3$ , bias = False, padding = same)
ReLU
Batch Normalization Layer (gain = 1, bias = 0)
Convolutional Layer (channels = 128, kernel-size = $3 \times 3$ , bias = False, padding = same)
ReLU
Max Pooling Layer (size = $2 \times 2$ , stride = $2 \times 2$ )
Batch Normalization Layer (gain = 1, bias = 0)
Convolutional Layer (channels = 256, kernel-size = $3 \times 3$ , bias = False, padding = same)
ReLU
Batch Normalization Layer (gain = 1, bias = 0)
Convolutional Layer (channels = 256, kernel-size = $3 \times 3$ , bias = False, padding = same)
ReLU
Max Pooling Layer (size = $2 \times 2$ , stride = $2 \times 2$ )
Batch Normalization Layer (gain = 1, bias = 0)
Convolutional Layer (channels = 512, kernel-size = $3 \times 3$ , bias = False, padding = same)
ReLU
Batch Normalization Layer (gain = 1, bias = 0)
Convolutional Layer (channels = 512, kernel-size = $3 \times 3$ , bias = False, padding = same)
ReLU
Max Pooling Layer (size = $2 \times 2$ , stride = $2 \times 2$ )
Batch Normalization Layer (gain = 1, bias = 0)
Fully Connected Layer (units = 1024, bias = False)
ReLU
Batch Normalization Layer (gain = 1, bias = 0)
Fully Connected Layer (units = 1024, bias = False)
ReLU
Batch Normalization Layer (gain = 1, bias = 0)
Fully Connected Layer (units = 1024, bias = False)
Batch Normalization Layer (gain = 1, bias = 0)
Softmax

Table 9. Algorithm specific train settings for MNIST, CIFAR-10, and CIFAR-100.

Algorithm	Setting	MNIST	CIFAR-10	CIFAR-100
BayesBiNN	Learning rate start $\alpha_0$	$10^{-4}$	$3 \cdot 10^{-4}$	$3 \cdot 10^{-4}$
	Learning rate end $\alpha_T$	$10^{-16}$	$10^{-16}$	$10^{-16}$
	Learning rate decay	Cosine	Cosine	Cosine
	MC train samples $S$	1	1	1
	MC test samples $C$	0	0	0
	Temperature $\tau$	$10^{-10}$	$10^{-10}$	$10^{-8}$
	Prior $\lambda_0$	<b>0</b>	<b>0</b>	<b>0</b>
	Initialization $\lambda$	$\pm 10$ randomly	$\pm 10$ randomly	$\pm 10$ randomly
STE Adam	Learning rate start $\alpha_0$	$10^{-2}$	$10^{-2}$	$10^{-2}$
	Learning rate end $\alpha_T$	$10^{-16}$	$10^{-16}$	$10^{-16}$
	Learning rate decay	Cosine	Cosine	Cosine
	Gradient clipping	Yes	Yes	Yes
	Weights clipping	Yes	Yes	Yes
Bop	Threshold $\tau$	$10^{-8}$	$10^{-8}$	$10^{-9}$
	Adaptivity rate $\gamma$	$10^{-5}$	$10^{-4}$	$10^{-4}$
	$\gamma$ -decay type	Step	Step	Step
	$\gamma$ -decay rate	$10^{\frac{-3}{500}}$	0.1	0.1
	$\gamma$ -decay interval (epochs)	1	100	100
PMF	Learning rate start	$10^{-3}$	$10^{-2}$	$10^{-2}$
	Learning rate decay type	Step	Step	Step
	LR decay interval (iterations)	7k	30k	30k
	LR-scale	0.2	0.2	0.2
	Optimizer	Adam	Adam	Adam
	Weight decay	0	$10^{-4}$	$10^{-4}$
	$\rho$	1.2	1.05	1.05
Adam (Full-precision)	Learning rate start $\alpha_0$	$3 \cdot 10^{-4}$	$10^{-2}$	$3 \cdot 10^{-3}$
	Learning rate end $\alpha_T$	$10^{-16}$	$10^{-16}$	$10^{-16}$
	Learning rate decay	Cosine	Cosine	Cosine

Table 10. Detailed results of different optimizers trained on MNIST, CIFAR-10 and CIFAR-100 (Averaged over 5 runs).

Dataset	Optimizer	Train Accuracy	Validation Accuracy	Test Accuracy
MNIST	STE Adam	99.78 ± 0.10 %	99.02 ± 0.11 %	<b>98.85 ± 0.09 %</b>
	Bop	99.23 ± 0.04 %	98.55 ± 0.05 %	98.47 ± 0.02 %
	PMF		99.06 ± 0.01 %	98.80 ± 0.06 %
	<b>BayesBiNN (mode)</b>	99.85 ± 0.05 %	99.02 ± 0.13 %	<b>98.86 ± 0.05 %</b>
	<b>BayesBiNN (mean)</b>	99.85 ± 0.05 %	99.02 ± 0.13 %	<b>98.86 ± 0.05 %</b>
	Full-precision	99.96 ± 0.02 %	99.15 ± 0.14 %	99.01 ± 0.06 %
CIFAR-10	STE Adam	99.99 ± 0.01 %	94.25 ± 0.42 %	<b>93.55 ± 0.15 %</b>
	Bop	99.79 ± 0.03 %	93.49 ± 0.17 %	93.00 ± 0.11 %
	PMF		91.87 ± 0.10 %	91.43 ± 0.14 %
	<b>BayesBiNN (mode)</b>	99.96 ± 0.01 %	94.23 ± 0.41 %	<b>93.72 ± 0.16 %</b>
	<b>BayesBiNN (mean)</b>	99.96 ± 0.01 %	94.23 ± 0.41 %	<b>93.72 ± 0.15 %</b>
	Full-precision	100.00 ± 0.00 %	94.54 ± 0.29 %	93.90 ± 0.17 %
CIFAR-100	STE Adam	99.06 ± 0.15 %	74.09 ± 0.15 %	72.89 ± 0.21 %
	Bop	90.09 ± 0.57 %	69.97 ± 0.29 %	69.58 ± 0.15 %
	PMF		69.86 ± 0.08 %	70.45 ± 0.25 %
	<b>BayesBiNN (mode)</b>	98.02 ± 0.18 %	74.76 ± 0.41 %	<b>73.68 ± 0.31 %</b>
	<b>BayesBiNN (mean)</b>	98.02 ± 0.18 %	74.76 ± 0.41 %	<b>73.65 ± 0.41 %</b>
	Full-precision	99.89 ± 0.02 %	75.89 ± 0.41 %	74.83 ± 0.26 %

Table 11. Test accuracy of BayesBiNN and LR-net trained on MNIST, CIFAR-10. Results for BayesBiNN are averaged over 5 random runs.

Optimizer	MNIST	CIFAR-10
LR-net Shayer et al. (2018)	99.47 %	93.18%
<b>BayesBiNN (mode)</b>	99.50 ± 0.02 %	93.97 ± 0.11 %

Table 12. The MLP architecture used for continual learning (Nguyen et al., 2018)

Fully Connected Layer (units = 100, bias = False)
ReLU
Batch Normalization Layer (gain = 1, bias = 0)
Fully Connected Layer (units = 100, bias = False)
ReLU
Batch Normalization Layer (gain = 1, bias = 0)
Fully Connected Layer (units = 100, bias = False)
ReLU
Batch Normalization Layer (gain = 1, bias = 0)
Softmax

Table 13. Algorithm specific train settings for continual learning on permuted MNIST.

Algorithm	Setting	Permuted MNIST
BayesBiNN	Learning rate start $\alpha_0$	$10^{-3}$
	Learning rate end $\alpha_T$	$10^{-16}$
	Learning rate decay	Cosine
	MC train samples $S$	1
	MC test samples $C$	100
	Temperature $\tau$	$10^{-2}$
	Prior $\lambda_0$	learned $\lambda$ of the previous task
	Initialization $\lambda$	$\pm 10$ randomly
	Batch size	100
	Number of epochs	100