
[Re:] Training Binary Neural Networks using the Bayesian Learning Rule

Anonymous Author(s)

Affiliation

Address

email

Reproducibility Summary

1

2 *(1) gives a mathematically principled approach to solve the discrete optimization problem that occurs in the case of*
3 *Binary Neural Networks and claims to give a similar performance on various classification benchmarks such as MNIST,*
4 *CIFAR-10, and CIFAR-100 as compared to their full-precision counterparts, as well as other recent algorithms to*
5 *train BNNs like PMF and Bop. The paper also claims that the BayesBiNN method has an application in the continual*
6 *learning domain as it helps in overcoming catastrophic forgetting of the past by using the posterior approximation of*
7 *the previous task as a prior for the upcoming task. We try to reproduce all the results presented in the original paper by*
8 *making a separate and independent codebase.*

9 **Scope of Reproducibility**

10 We try to verify the performance of our re-implementation of the BayesBiNN optimizer on various classification and
11 regression benchmarks. We also implemented the STE optimizer which was the central baseline model used in the
12 paper. Finally, we tried to evaluate the results of BayesBiNN on the continual learning benchmark to get a better insight.

13 **Methodology**

14 We developed our separate code-base, consisting of an end-to-end trainer with a Keras-like interface, for the reproduction
15 which includes the implementation of the BayesBiNN and STE optimizer. We did refer to the author's code open-sourced
16 on GitHub to get some insights about the hyperparameters and other doubts that emerged during code development.

17 **Results**

18 We reproduced the accuracy of the BayesBiNN optimizer within less than 0.5% of the originally reported value, which
19 upholds the conclusion that it performs nearly as well as its full-precision counterpart in classification tasks. When we
20 tried this in a semantic segmentation context, we found that the results were very underwhelming and in contrast with
21 the seemingly good results by the STE optimizer even with much hyperparameter tuning. We can conclude that, like
22 other Bayesian methods, it is difficult to train BayesBiNN on more complex tasks.

23 **What was easy**

24 After we worked out the mathematics behind the BayesBiNN approach, we developed a pseudo-code for the optimization
25 process which along with references from the author's code, helped us a lot in our reproduction study.

26 **What was difficult**

27 Some of the hyperparameters were not mentioned by the authors in their paper so it was difficult to approximate the
28 values of those parameters. The lack of resources was the next big difficulty that we faced.

29 **Communication with original authors**

30 We had a very fruitful conversation with the authors, which helped us in better understanding the BayesBiNN approach
31 and its extension to the segmentation domain. The detailed pointers are given at the end of this report.

32 1 Introduction

33 Deep Learning is moving towards larger and larger parameters day-by-day, which often makes it difficult to run on
34 resource-constraint devices like mobile phones. Binary Neural Networks (BNNs) could act as a savior in such situations,
35 helping in largely saving storage and computational costs. The problem of optimizing this binary set of weights
36 is clearly a discrete optimization problem. Previous approaches like Straight-Through Estimator (STE) and Binary
37 Optimizer (Bop) tend to ignore this and use gradient-based methods, which still worked in practice. The paper presents
38 a mathematically principled approach for training BNNs which also justifies the current approaches.

39 2 Scope of reproducibility

40 The paper mentions a bayesian approach to solve the discrete optimization problem in the case of Binary Neural
41 Networks (BNNs). The outcome of this approach was a BayesBiNN optimizer which could be used to train BNNs and
42 achieve similar accuracy as compared to their full-precision counterparts. To verify the claims given in the paper, we
43 target to achieve the following objectives:

- 44 • Work out and present the mathematics behind BayesBiNN in a simpler way and present the pseudo-code to
45 the optimizer.
- 46 • Implement the BayesBiNN optimizer and STE optimizer to verify the accuracy on tasks of varying complexities,
47 as reported in the original paper.
- 48 • Reproduce the results for other baselines present in the paper such as proximal mean-field (PMF) according to
49 the hyper-parameters given in the paper.
- 50 • Evaluating the performance of BayesBiNN optimizer in more complex domains like semantic segmentation.

51 3 Methodology

52 We have re-implemented the algorithm proposed in the paper from scratch using PyTorch and created an end-to-
53 end model trainer with a Keras-like interface. We referred to the code given by the authors for the baseline model
54 hyperparameters and the source of synthetic datasets. Following is the algorithmic form of what the authors have
55 presented in the paper.

Algorithm 1: Bayesian Learning rule for BayesBiNN

Input: Initialize λ

for *number of training epochs* **do**

for $i = 1, \dots, \text{number of mini-batch examples}$ **do**

 Sample $\epsilon \sim \mathcal{U}(0, 1)$ and set $\delta = \frac{1}{2} \log \frac{\epsilon}{1-\epsilon}$

 Initialize $w_b = \tanh((\lambda + \delta)/\tau)$

 Compute following using *gumbel-softmax* trick

$$g_i := \frac{1}{M} \nabla_{w_b} l(y_i, f_{w_r}(x_i))$$

$$s_i := \frac{N(1 - w_b^2)}{\tau(1 - \tanh(\lambda)^2)}$$

end

 Update μ and λ using following equation

$$\mu \leftarrow \tanh(\lambda)$$

$$\lambda \leftarrow (1 - \alpha)\lambda - \alpha \left[\sum_{i=1}^M (s_i \odot g_i) - \lambda_0 \right]$$

end

57 This would make the paper more interpretable in terms of implementation.

58 Some of the mathematical expressions mentioned in the original paper were presented from various sources and missed
59 out several intermediate steps which we found to be very important while reproducing the paper from scratch. Here we
60 present a step-wise derivation of some important expressions written in the original paper:

61 Bayesian formulation of the discrete optimization problem, in which loss has to be minimized w.r.t posterior $q(w)$,
 62 given prior $p(w)$ can be written as:

$$\mathbb{E}_{q(w)}\left[\sum_{i=1}^N l(y_i, f_w(x_i))\right] + \mathcal{D}_{KL}[q(w)||p(w)]$$

63 To solve the above optimization problem, Bayesian learning rule given in (6) is applied, assuming solution to be a part
 64 of minimal exponential family of distribution, given by:

$$q(w) = h(w)\exp[\lambda^T \phi(w) - A(\lambda)]$$

65 where base measure $h(w)$ is assumed to be 1. Following is the update rule used to learn λ :

$$\lambda \leftarrow (1 - \rho)\lambda - \rho[\nabla_{\mu} \mathbf{E}_{q(w)}[l(y_i, f_w(x_i))] - \lambda_0]$$

66 where ρ is the learning rate, $\mu = \mathbf{E}_{q(w)}[\phi(w)]$. Bernoulli distribution being a special case of minimal exponential
 67 family distribution, we assume prior $p(w) \sim \text{Bern}(p)$ with $p = 0.5$, and posterior $q(w)$ to be mean-field Bernoulli
 68 distribution:

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

69 For weight j ,

$$\begin{aligned} q(w_j) &= \exp\left(\frac{1}{2}(1 + w_j) \log p_j + \frac{1}{2}(1 - w_j) \log(1 - p_j)\right) \\ &= \exp\left(\underbrace{w_j}_{\phi(w)} \underbrace{\frac{1}{2} \log \frac{p}{1-p}}_{\lambda} + \frac{1}{2} \log(p(1-p))\right) \end{aligned}$$

71 Comparing above expression with minimal exponential family distribution, we can say:

$$\lambda = \frac{1}{2} \log \frac{p}{1-p} \text{ and } \phi(w) = w.$$

72 We defined $\mu = \mathbb{E}_{q(w)}[\phi(w)]$,

$$\begin{aligned} \mu &= \int wq(w)dw = \mathbb{E}[q(w)] = \sum_{w^i \in \{-1,1\}} w^i q(w^i) \\ &= \sum_{w^i \in \{-1,1\}} w^i p^{\frac{1+w^i}{2}} (1-p)^{\frac{1-w^i}{2}} = -(1-p) + p \\ &= 2p - 1 \end{aligned}$$

75 From above derivations we can say that, $p = 1/(1 + \exp(-2\lambda)) = \text{Sigmoid}(2\lambda)$ and $q(w) \sim \text{Bern}(p)$.

76 To implement the update rule, we need to compute the gradient with respect to μ . Original paper uses a reparamateriza-
 77 tion trick called Gumbel-softmax trick (7), which is used to relax the discrete random variables of a concrete distribution
 78 (for eg, Bernoulli distribution). Binary concrete relaxation (7) of Binary concrete random variable $X \in (0, 1)$ with
 79 distribution $X \sim \text{BinConcrete}(\alpha, \lambda)$ with temperature λ and location α ,

$$X = \frac{1}{1 + \exp(-(\log \alpha + L)/\lambda)}$$

80 where $L \sim \text{Logistic}$. And its density is given by

$$p_{\alpha, \lambda}(x) = \frac{\lambda \alpha x^{-\lambda-1} (1-x)^{-\lambda-1}}{(\alpha x^{-\lambda} + (1-x)^{-\lambda})^2}$$

81 Using above expressions, for binary weights $w_j \in \{0, 1\}$, relaxed variable $w_r^{\epsilon_j, \tau}(p_j) \in (0, 1)$ can be used with
 82 temperature τ and $\alpha = e^{2\lambda}$ given by

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

83 where $\delta_j \sim \text{Logistic}$ and its density is given by

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

84 4 Experimental setup

85 4.1 Model descriptions

86 We kept the model architectures the same as mentioned in the original paper to maintain uniformity and implemented
87 them ourselves. For the MNIST classification task, we used the BinaryConnect architecture and for the CIFAR
88 classification task, we use the VGGBinaryConnect architecture. The authors also compared their BayesBiNN method
89 with the LR-Net method in (8). We implemented the same model architecture as in the LR-Net paper. The detailed
90 architectures are mentioned in the supplementary material provided with this report. For the segmentation task, we used
91 the original U-Net architecture, mentioned in (11) with a minor difference that we introduced a BatchNorm layer after
92 every convolution layer.

93 4.2 Datasets

94 The datasets used for image classification tasks are MNIST, CIFAR-10, and CIFAR-100. For generating visualizations
95 for the BayesBiNN and STE methods, we used small toy datasets, the Snelson dataset (10) for regression problems, and
96 Two Moon’s dataset (9) for classification problems. For the segmentation part, we used the Brain Tissue segmentation
97 dataset taken from (11), and for the continual learning visualizations, we used the permuted MNIST dataset (12). The
98 pre-processing of inputs has been kept the same as mentioned in the original paper and has been detailed below.

99 **Pre-processing:** For the MNIST dataset we simply normalize the images and do not perform data augmentation.
100 We keep our validation split as 0.1 uniformly across all sets of experiments except the comparison with the LR-Net
101 method (8). For the CIFAR datasets also, we perform the normalization of images along with data-augmentation where
102 we generate images by randomly cropping a 32x32 image from a 40x40 padded image. Finally, for our semantic
103 segmentation task, we had a very small dataset of 30 images, out of which 24 were chosen for training and 6 for
104 obtaining the validation score. No other pre-processing has been done in this case.

105 4.3 Hyperparameters

106 We have used the hyper-parameters given in the original paper. Table 1 contains the list of all the parameters we used
107 for our experiments:

Optimizer	Parameter	MNIST	CIFAR10	CIFAR100	Snelson Dataset	2 Moons Dataset
BayesBiNN	MC steps	1	1	1	1	5
	Initial LR	10^{-4}	3.10^{-4}	3.10^{-4}	10^{-4}	10^{-3}
	Final LR	10^{-16}	10^{-16}	10^{-16}	10^{-5}	10^{-5}
	LR Scheduler	Cosine	Cosine	Cosine	MultiStepLR	MultiStepLR
	Temperature τ	10^{-10}	10^{-10}	10^{-8}	1	1
	Initialization λ	± 10	± 10	± 10	± 10	± 15
STE	Initial LR	10^{-2}	10^{-2}	10^{-2}	10^{-1}	10^{-1}
	Final LR	10^{-16}	10^{-16}	10^{-16}	10^{-1}	10^{-3}
	LR Scheduler	Cosine	Cosine	Cosine	MultiStepLR	MultiStepLR
Adam (Full Precision)	Initial LR	10^{-5}	10^{-4}	10^{-4}	-	-
	Final LR	Step	Step	Step	-	-
	LR Scheduler	1	100	100	-	-

Table 1: Training setting for different optimizers on MNIST, CIFAR10, and CIFAR100 datasets.

108 4.4 Computational requirements

109 All our final experimental results were performed on a machine having 1 NVIDIA Tesla V100 GPU and 1 single-core
110 system with 16 GB memory. Training the Binary Network with BayesBiNN optimizer for a single run, takes around 2.5
111 GPU hours for MNIST, 5.5 GPU hours for CIFAR-10, and around 8.5 GPU hours for the CIFAR-100 dataset, in the
112 current experimental setup.

113 5 Results

114 In Table 1 we report our results for various classification benchmarks using our implemented BayesBiNN and STE
115 optimizer. We notice that we get a difference of less than 0.1% as compared to that in the original paper. We generated

116 the results for baseline STE optimizer and full-precision networks by evaluating our implementation of these methods.
 117 We also generated the results of PMF, by modifying its original open-sourced code and using the hyperparameters
 118 mentioned in the original paper.

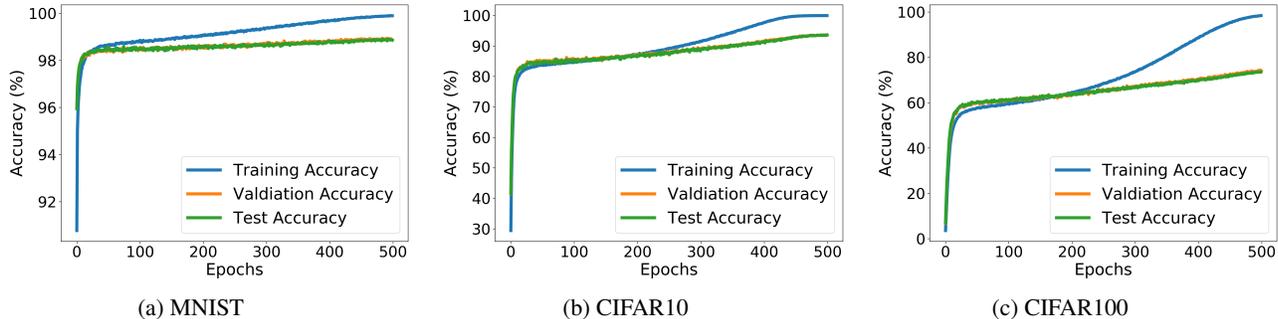


Figure 1: Training/Validation/Test accuracy using BayesBiNN optimizer

Datasets	Optimizer	Training Accuracy	Validation Accuracy	Test Accuracy
MNIST	BayesBiNN(ours)	99.90 ±0.01%	99.89 ±0.07%	98.87 ±0.06%
	BayesBiNN(orig.)	99.85 ±0.05%	99.02 ±0.13%	98.86 ±0.05%
	STE	99.90 ±0.01%	98.86 ±0.09%	98.89 ±0.05%
	PMF	-	98.73%	-
	Adam (Full Precision)	99.98 ±0.01%	99.02 ±0.04%	99.02 ±0.01%
CIFAR10	BayesBiNN(ours)	99.96 ±0.01%	93.59 ±0.45%	93.54 ±0.26%
	BayesBiNN(orig.)	99.96 ±0.01%	94.23 ±0.41%	93.72 ±0.16%
	STE	99.99 ±0.01%	93.77 ±0.06%	93.54 ±0.08%
	PMF	-	91.98%	-
	Adam (Full Precision)	99.99 ±0.01%	94.27 ±0.15%	94.38 ±0.16%
CIFAR100	BayesBiNN(ours)	98.35 ±0.1%	74.13 ±0.78%	73.56 ±0.06%
	BayesBiNN(orig.)	98.02 ±0.18%	74.76 ±0.41%	73.68 ±0.31%
	STE	99.22 ±0.03%	72.74 ±0.06%	73.25 ±0.26%
	PMF	-	70.82%	-
	Adam (Full Precision)	99.89 ±0.02%	75.04 ±0.71%	74.80 ±0.39%

Table 2: Results of different optimizers trained on MNIST, CIFAR10, and CIFAR100.

119 5.1 Comparison with LR-Net

120 Authors compared their BayesBiNN approach to the LR-Net method presented in (8). We tried to reproduce the
 121 result for the same setting. In this comparison, the data pre-processing and augmentation methods remain the same as
 122 mentioned in section 4.2, but we do not split the data in training and validation sets in this case. We denote the test
 123 accuracies after 190 epochs in the case of MNIST and 290 epochs in the case of CIFAR-10, as done in the original
 124 paper to maintain uniformity. Note that, our accuracy is matching with that of the original authors in the case of MNIST
 125 but not in the case of CIFAR-10. We suspect that this is due to some difference in Batch-Norm layers used.

Optimizer	MNIST	CIFAR10
BayesBiNN (ours)	99.52%	84.49%
BayesBiNN (orig.)	99.50%	93.97%
LR-net (8)	99.47%	93.18%

Table 3: Test accuracy of BayesBiNN and LRNet.

126 5.2 Continual Learning

127 As mentioned in the original paper, we try to reproduce the author’s claims about weight distribution across tasks in a
 128 simple continual learning domain tested on Permuted MNIST. As we can see, as we learn across the tasks, the curve

129 becomes flat from the middle conveying that the weights become more deterministic. Our result matches with the
 130 claims in the original paper.

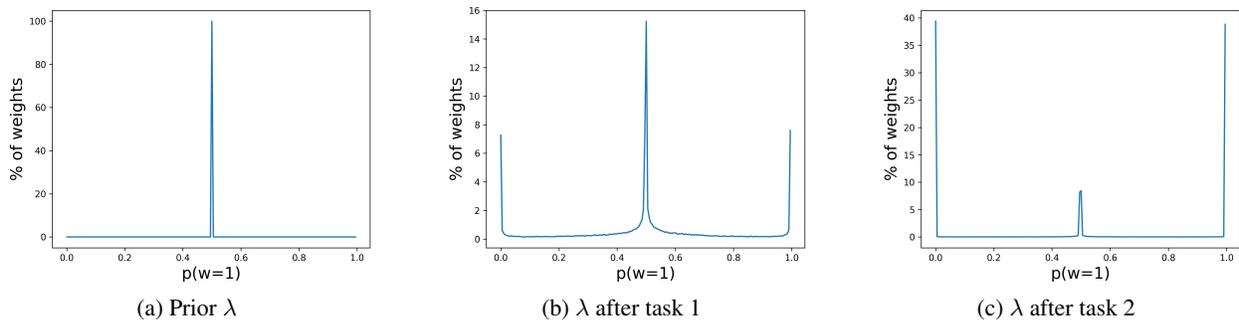


Figure 2: Distribution of $p(w = 1)$ across consecutive learning tasks

130

131 5.3 Visualization using Synthetic Dataset

132 In the original paper, the authors present visualizations on binary classification (Two moons dataset(9)) and toy
 133 regression (Snelson dataset(10)) using STE and BayesBiNN optimizer. For the classification task, the authors claimed
 134 that STE is a more deterministic classifier compared to BayesBiNN. We reproduced this experiment and the results
 135 depicted in Figure 3 seem to be consistent with the author’s claim. For the regression task, we conclude that the author’s
 136 claim about BayesBiNN (mean) giving a smoother curve compared to STE is true, which can also be seen in Figure 4.

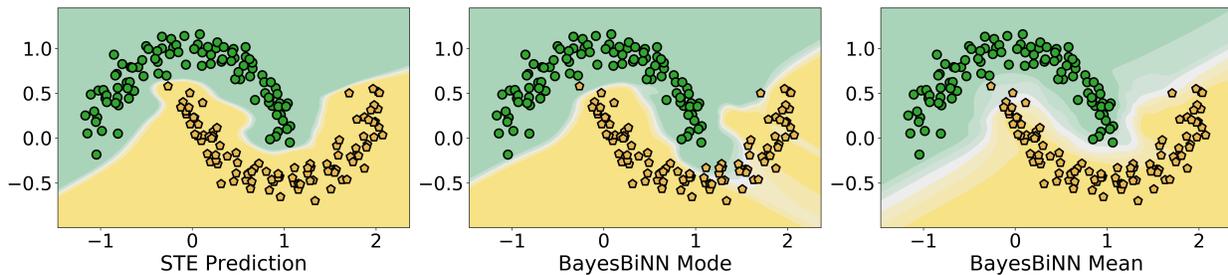


Figure 3: Classification on Two Moons dataset using STE and BayesBiNN optimizer.

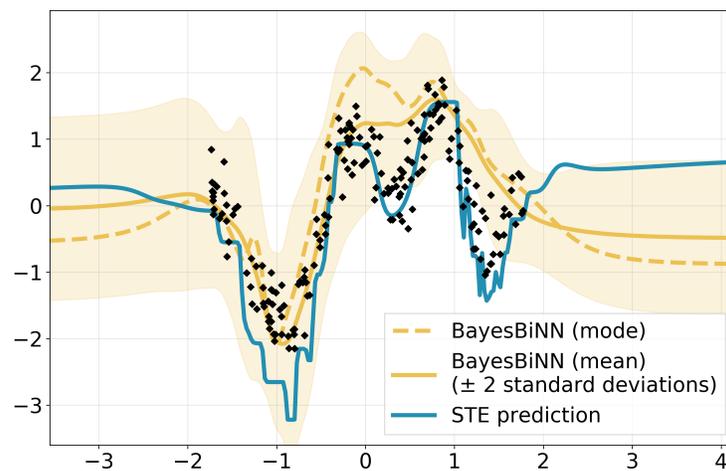


Figure 4: Regression on Snelson dataset using STE and BayesBiNN optimizer.

Temperature	10	1	0.1	10^{-2}	10^{-3}	10^{-4}
MSE Loss	1.313	0.208	2.151	0.443	0.231	0.199
Temperature	10^{-5}	10^{-6}	10^{-7}	10^{-8}	10^{-9}	10^{-10}
MSE Loss	0.156	0.127	0.173	0.122	0.195	0.173

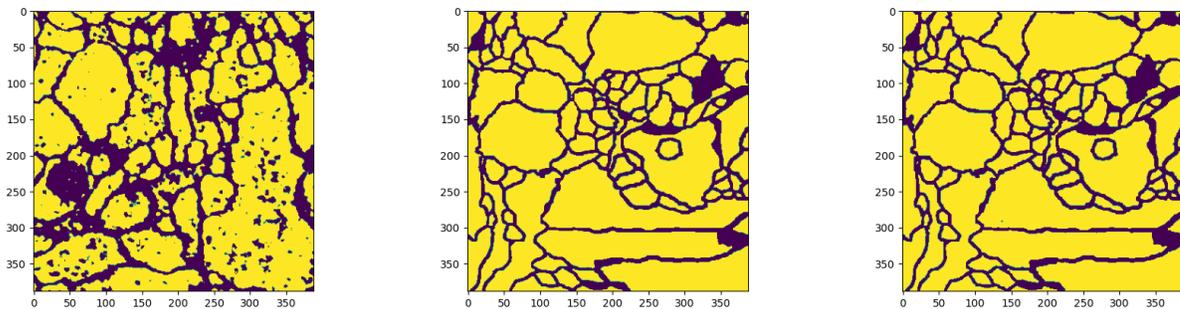
Table 4: Mean square error loss of Snelson dataset for different temperatures.

137 5.4 Extended Results (Semantic Segmentation)

138 We tried to validate the performance of the BayesBiNN optimizer on more complex tasks like Semantic Segmentation.
 139 Unfortunately, the results with BayesBiNN were quite underwhelming as compared to STE and its full-precision
 140 counterpart. We tried various parameters to improve its performance but none seemed to work. We had a brief
 141 discussion with the authors regarding this issue and the authors suggested that Bayesian models are intrinsically very
 142 difficult to train. For the results shown in Table 5 and Figure 5, we have used the hyperparameters denoted in Table 1.

	BayesBiNN	STE	Adam (Full Precision)
Validation Score	0.4102	0.3108	0.2943

Table 5: (1 - IoU) score for validation set



(a) BayesBiNN

(b) STE

(c) Adam

Figure 5: Some samples of segmented image outputs

143 6 Discussion

144 We reproduced almost all the experiments given in the original paper and most of our results match with the original
 145 claims. While this BayesBiNN approach is mathematically principled, we tried to take a step forward, by using
 146 that optimizer on a single segmentation task. However, the results were against our expectation and the result of
 147 segmentation was a zoomed segmented image of the input with lots of noise. Apart from this, even in the case of
 148 comparison with the LR-Net method, our accuracy differs from that of the original authors, which we feel might be due
 149 to some difference in architecture chosen. The major contribution of our work is developing a code base library based
 150 on PyTorch with a Keras type interface for training BNNs with several different methods in its arsenal. This would
 151 reduce the coding efforts while training a BNNs and could help in future research as benchmarking platform.

152 6.1 What was easy

153 The original paper contained a very good explanation of the mathematics behind the BayesBiNN approach. After we
 154 worked that out the pseudo-code as pointed out in Algorithm 1, the basic implementation of the optimizer became
 155 easy and easily verifiable by the author’s original code. The appendix in the original paper contained a list of various
 156 hyper-parameters used for experiments. This helped us a lot while running the experiments and deciding the range of
 157 hyper-parameters while doing ablation studies.

158 6.2 What was difficult

159 The most difficult part here was running a large number of experiments in lack of many computational resources. This
160 difficulty was increased since we are taking an average of 5 runs while reporting all our results. Apart from this, we
161 also faced some difficulty in taking care of the hyper-parameters, which were not mentioned in the original paper
162 (like momentum coefficient). To cater to that, we had to guess some possible values of the hyper-parameters and run
163 small random searches to find a good candidate. Finally, we also faced difficulty while reproducing the results for the
164 baselines PMF and Bop and adapting their experimental settings to match with those used in the original BayesBiNN
165 paper. Since their code was written a long time ago and used older technologies, this task took us a lot of time.

166 6.3 Communication with original authors

167 We did not understand the intent of the authors for choosing temperature as 1 in the case of experiments on synthetic
168 datasets. We were also curious about the author’s view on segmentation tasks using BayesBiNN. We mailed this, along
169 with the review of their paper, to the authors to ask for some pointers. They gave the following major pointers:

- 170 • It is reasonable that at high temperatures the learned distribution will have high variance. The mode mentioned
171 in the paper refers to the $\text{sign}(\hat{w})$, where \hat{w} denotes the expectation of the learned posterior Bernoulli
172 distribution. It is not appropriate to directly use the continuous \hat{w} as the mode. Another way is to use mean,
173 which samples from the learned posterior Bernoulli distribution, and then make predictions using ensemble
174 learning.
- 175 • STE is more stable and suggested by the authors to act as a baseline, in particular, Adam STE first, to make
176 sure binary networks work. As shown in the paper, there is literally very little difference between STE and
177 BayesBiNN but indeed the latter is difficult to train, as most Bayesian optimizers.

178 Broader Impact

179 Recent researches (3) mention that training a single big transformer model could emit around 626,155 lbs CO₂ which is
180 around 5 times of average carbon emission by a car in its total lifetime. Clearly, Deep Learning takes a huge toll on
181 the environment which is why there has been an increased focus on much more energy-efficient "Green AI". BNNs
182 intrinsically have far less computational and space complexity as compared to their full-precision counterparts and
183 as we can see above they can also achieve accuracy close to the full-precision networks, at least in the classification
184 tasks, and also show the potential of expanding well to more complex segmentation tasks. This can help us a lot in
185 moving towards cleaner Deep Learning. This class of technology also provides a huge set of opportunities in extending
186 AI to edge devices with much smaller and low-energy systems. We feel that its potential impact on the environment
187 and sustainability is at par with its academic importance, that is why we see it as a much larger thing than just a set of
188 publications.

189 References

- 190 [1] Xiangming Meng, Roman Bachmann & Mohammad Emtiyaz Khan, Training Binary Neural Networks using the
191 Bayesian Learning Rule. In *Proceedings of the 37th International Conference on Machine Learning*, Online,
192 PMLR 119, 2020.
- 193 [2] L. Biewald, “Experiment Tracking with Weights and Biases,” *Weights & Biases*. [Online]. Available:
194 <http://wandb.com/>.
- 195 [3] Emma Strubell, Ananya Ganesh & Andrew McCallum, Energy and Policy Considerations for Deep Learning in
196 NLP. In the *57th Annual Meeting of the Association for Computational Linguistics (ACL)*. Florence, Italy. July 2019
- 197 [4] Ajanthan, T., Dokania, P. K., Hartley, R., & Torr, P. H. S, Proximal mean-field for neural network quantization. In
198 *IEEE International Conference on Computer Vision*, pp. 4871–4880, 2019a.
- 199 [5] Helwegen, K., Widdicombe, J., Geiger, L., Liu, Z., Cheng, K.-T., & Nusselder, R., Latent weights do not exist:
200 Rethinking binarized neural network optimization. In *NeurIPS*, 2019.
- 201 [6] Khan, M. E. and Rue, H., Learning-Algorithms from Bayesian principles 2020. [https://emtiyaz.github.io/
202 papers/learning_from_bayes.pdf](https://emtiyaz.github.io/papers/learning_from_bayes.pdf).

- 203 [7] Chris J. Maddison, Andriy Mnih, Yee Whye Teh, The Concrete Distribution: A Continuous Relaxation of Discrete
204 Random Variables In *ICLR*, 2017
- 205 [8] Shayer, O., Levi, D., and Fetaya, E. Learning discrete weights using the local reparameterization trick. *ICLR*, 2018.
- 206 [9] Moons, T. Two moons datasets description. [https://scikit-learn.org/stable/modules/generated/
207 sklearn.datasets.make_moons.html](https://scikit-learn.org/stable/modules/generated/sklearn.datasets.make_moons.html).
- 208 [10] Snelson, E. and Ghahramani, Z. Sparse Gaussian processes using pseudo-inputs. In *Proceedings of the 18th
209 International Conference on Neural Information Processing Systems*, NIPS05, pp. 1257-1264, Cambridge, MA,
210 USA, 2005. MIT Press.
- 211 [11] Ronneberger, Olaf, Philipp Fischer, and Thomas Brox. "U-net: Convolutional networks for biomedical image
212 segmentation." *International Conference on Medical image computing and computer-assisted intervention*. Springer,
213 Cham, 2015.
- 214 [12] Goodfellow, I. J., Mirza, M., Xiao, D., Courville, A., and Bengio, Y. An empirical investigation of catastrophic
215 forgetting in gradient-based neural networks. *ICLR*, 2013.
- 216 [13] Zenke, F., Poole, B., and Ganguli, S. Continual learning through synaptic intelligence. In Proceedings of the 34th
217 International Conference on Machine Learning-Volume 70, pp. 3987-3995. JMLR. org, 2017.
- 218 [14] LeCun, Y. and Cortes, C. MNIST handwritten digit database. 2010. URL [http://yann.lecun.com/exdb/
219 mnist/](http://yann.lecun.com/exdb/mnist/).
- 220 [15] Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *ICLR*, 2015.
- 221 [16] Jang, E., Gu, S., and Poole, B. Categorical reparameterization with Gumbel-softmax. *ICLR*, 2017.
- 222 [17] Courbariaux, M., Bengio, Y., and David, J.-P. Binaryconnect: Training deep neural networks with binary weights
223 during propagations. In *Advances in neural information processing systems*, pp. 3123-3131, 2015.