

Improving Bayesian Neural Networks by Adversarial Sampling

Jiaru Zhang¹ Yang Hua² Tao Song¹ Hao Wang³ Zhengui Xue¹ Ruhui Ma¹
Haibing Guan¹

¹Shanghai Jiao Tong University ²Queen's University Belfast

³Louisiana State University

- Bayesian neural networks have shown considerable potential and has been widely used in many tasks.
- Bayesian neural networks have indeed few publicized deployments in industrial practice despite the theoretical advancements.
- It is still unknown that why Bayesian neural networks can not learn a suitable representation and perform well.
- In this paper, we present a reason and propose Adversarial Sampling as a solution.

- The learning target of Bayesian neural networks with variational inference is

$$\begin{aligned}\mathcal{L} &= - \int Q_{\theta}(\mathbf{W}) \log \frac{P(\mathbf{W}, \mathcal{D})}{Q_{\theta}(\mathbf{W})} d\mathbf{W} \\ &= \underbrace{-\mathbb{E}_{\mathbf{W} \sim Q_{\theta}(\mathbf{W})} \log P(\mathcal{D} | \mathbf{W})}_{\mathcal{L}_p} + \underbrace{KL(P(\mathbf{W}) || Q_{\theta}(\mathbf{W}))}_{\mathcal{L}_r}.\end{aligned}$$

- It can be divided into two terms.
- The first term \mathcal{L}_p is directly related to the predictions.
- The second term \mathcal{L}_r can be seen as a regularization on the model parameters.

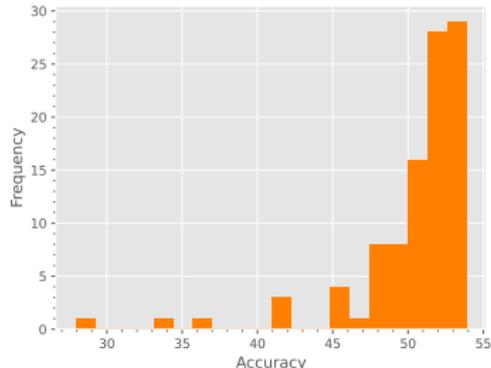
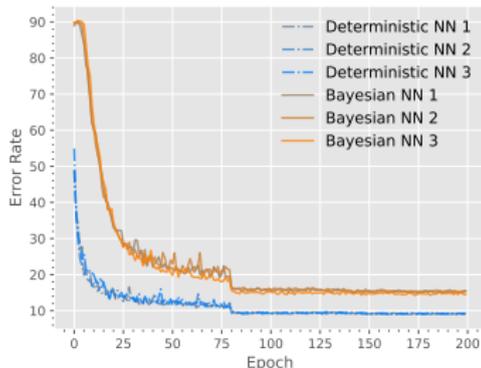
Explanation of the Poor Performance

Because of the randomness of sampling during training and testing,

- There are some errors in updating the parameters.
- Some models with poor performance are yielded in random sampling.

Validation:

- The curves of Bayesian neural networks fluctuate more sharply.
- Some models have much lower accuracies compared with others.



For dataset \mathcal{D} , parameter distributions $Q_\theta(\mathbf{W})$, we define

Adversarial Distribution

$$Q_{adv} = \operatorname{argmax}_{W[Q_{adv}, Q_\theta] \leq d} - \mathbb{E}_{W \sim Q_{adv}(W)} \log P(\mathcal{D} | \mathbf{W}). \quad (1)$$

- $W[Q_{adv}, Q_\theta]$ denotes the Wasserstein distance between Q_{adv} and Q_θ .
- d is a hyperparameter to control $W[Q_{adv}, Q_\theta]$.

Corresponding to the Adversarial Distribution, the adversarial loss is defined as

Adversarial Loss

$$\mathcal{L}_{adv} = -\mathbb{E}_{W \sim Q_{adv}(W)} \log P(\mathcal{D}|\mathbf{W}) \quad (2)$$

The total learning target is

Total Learning Target

$$\theta = \underset{\theta}{\operatorname{argmin}} ((1 - \lambda) \cdot \mathcal{L}_p + \lambda \cdot \mathcal{L}_{adv} + \mathcal{L}_r) \quad (3)$$

- λ controls the ratio of training with Adversarial Distribution.

- The total learning target is equivalent to the original one when $d = 0$ or $\lambda = 0$.
- Sampling from the Adversarial Distribution yields likely models with the worst performance.
- Parameters updating accordingly guarantees the performance of the regularly sampled models.

Adversarial Sampling as an Approximation

- The calculation of Q_{adv} analytically is difficult.
- We propose an iterative approach, Adversarial Sampling, as an approximation.

Adversarial Sampling

1. Sample w_{adv} from the parameter distribution $N(\mu, \sigma^2)$.
 2. Repeat multiple times:
 - $w_{adv} = w_{adv} + \alpha \cdot \sigma \cdot \text{sign}(\text{grad}(w_{adv}))$.
- Generating w_{adv} can be regarded as a sampling.
 - Many W_{adv} s create an approximation of Q_{adv} .

Implementation of Adversarial Sampling

- With the reparameterization trick, w is from $w = \mu + \epsilon \cdot \sigma, \epsilon \sim \mathcal{N}(0, 1)$.
- It makes implementation easier and training with gradient descent possible by updating ϵ directly.

Algorithm 1: Training with Adversarial Sampling

Input: Variational posterior parameters (μ, σ) , Batch data \mathcal{D}

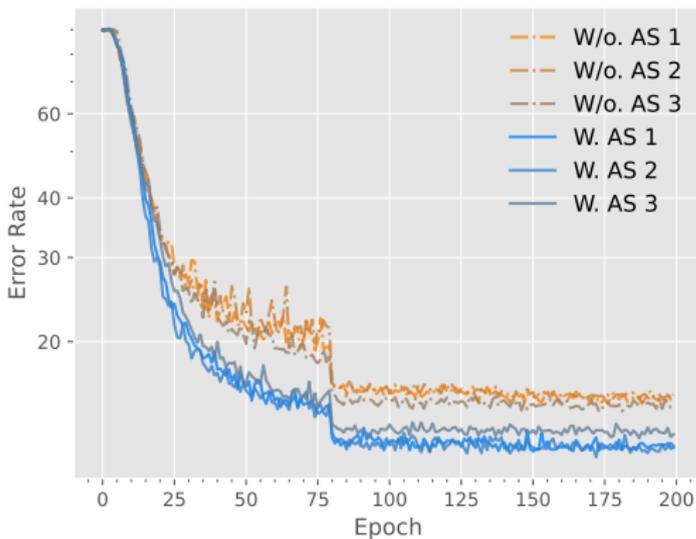
Parameters: Iterations N , Step length for perturbation α

Output: Updated variational posterior parameters (μ, σ)

- 1: Sample $\epsilon_{adv} \sim \mathcal{N}(0, 1)$
 - 2: **for** sufficient iterations N **do**
 - 3: Let $w_{adv} = \mu + \epsilon_{adv} \cdot \sigma$
 - 4: Calculate the adversarial loss \mathcal{L}_{adv} with parameter w_{adv} and data \mathcal{D}
 - 5: Update $\epsilon_{adv} = \epsilon_{adv} + \alpha \cdot \text{sign}(\frac{\partial \mathcal{L}_p}{\partial \epsilon_{adv}})$
 - 6: **end for**
 - 7: Let $w_{adv} = \mu + \epsilon_{adv} \cdot \sigma$
 - 8: Calculate the adversarial loss \mathcal{L}_{adv} with parameter w_{adv} and data \mathcal{D}
 - 9: Sample $\epsilon \sim \mathcal{N}(0, I)$
 - 10: Let $w = \mu + \epsilon \cdot \sigma$
 - 11: Calculate the prediction loss \mathcal{L}_p with parameter w and data \mathcal{D}
 - 12: Calculate the regularization loss \mathcal{L}_r analytically
 - 13: Calculate the total loss \mathcal{L} with Equation (11)
 - 14: Update parameter μ and σ with the total loss \mathcal{L}
-

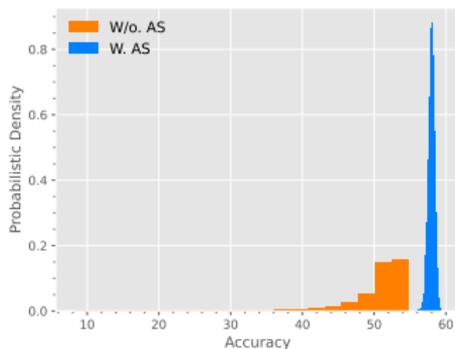
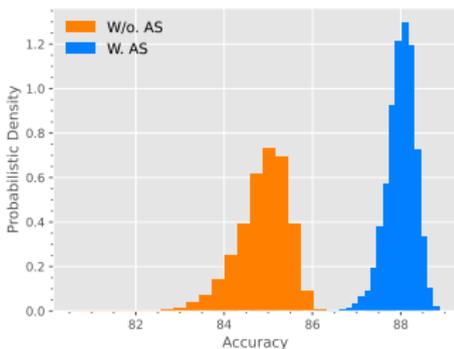
Experiments: Verification of Motivation

- Models trained with Adversarial Sampling have lower error rates.
- The change trends of models trained with Adversarial Sampling are more stable and steady.



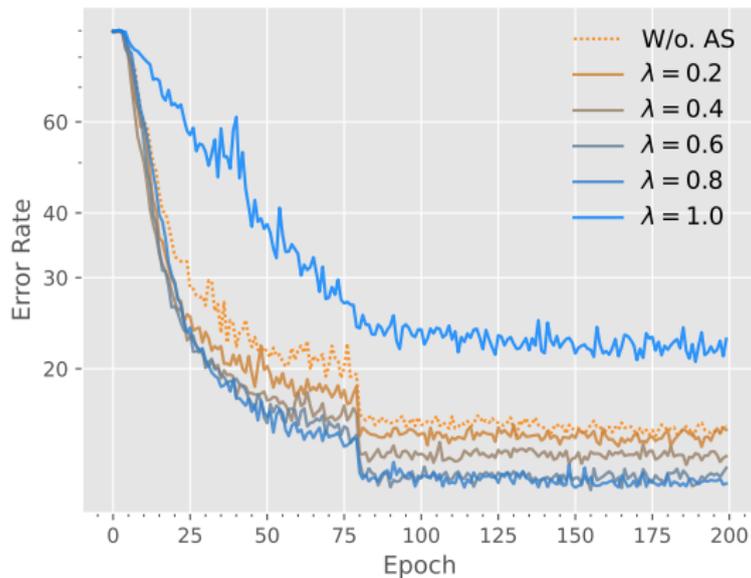
Experiments: Verification of Motivation

- The models trained without Adversarial Sampling distribute more dispersed.
- The accuracies of models trained with Adversarial Sampling are clearly higher.



Selection of Hyperparameter λ

- Model performance gets improved when λ increases from 0 to 0.8.
- There is a significant drop in performance when λ reaches 1.0.
- It validates the necessity of the introduction of parameter λ .



Improvement on Model Performance

We present three kinds of accuracies:

- The lowest accuracy and the highest accuracy among 100 sampled models.
- The accuracy of the ensembled model.

Models trained with Adversarial Sampling have much higher accuracies.

Dataset	Model	Lowest Accuracy	Highest Accuracy	Ensembled Accuracy
CIFAR-10	ResNet20	82.73 \pm 0.88	86.03 \pm 0.43	87.01 \pm 0.65
	ResNet20 + AS	86.33 \pm 0.45	88.35 \pm 0.51	88.76 \pm 0.73
	ResNet56	82.71 \pm 0.55	86.84 \pm 0.04	88.22 \pm 0.41
	ResNet56 + AS	87.30 \pm 0.32	88.86 \pm 0.79	89.61 \pm 0.93
	VGG	85.04 \pm 0.44	88.47 \pm 0.12	89.80 \pm 0.12
	VGG + AS	88.68 \pm 0.53	90.39 \pm 0.35	90.86 \pm 0.32
CIFAR-100	ResNet20	52.54 \pm 1.54	55.58 \pm 1.33	56.56 \pm 1.07
	ResNet20 + AS	54.83 \pm 0.95	57.24 \pm 0.89	57.62 \pm 0.91
	ResNet56	44.92 \pm 5.58	51.67 \pm 2.99	53.21 \pm 2.40
	ResNet56 + AS	54.76 \pm 2.26	57.50 \pm 1.43	58.63 \pm 1.58
	VGG	40.61 \pm 1.28	45.38 \pm 0.95	47.60 \pm 1.01
	VGG + AS	51.14 \pm 1.23	54.95 \pm 0.53	56.11 \pm 0.66

Combination with Bayesian Fine-tune

- Bayesian fine-tune is an effective method to improve the performance.
- Models trained with Adversarial Sampling also perform better under this higher baseline.

Dataset	Model	Lowest Accuracy	Highest Accuracy	Ensembled Accuracy
CIFAR-10	ResNet20	86.35 \pm 0.62	90.29 \pm 0.28	91.88 \pm 0.06
	ResNet20 + AS	88.19 \pm 0.44	91.22 \pm 0.18	91.98 \pm 0.18
	ResNet56	85.54 \pm 1.24	90.48 \pm 0.51	92.34 \pm 0.44
	ResNet56 + AS	88.74 \pm 0.64	91.78 \pm 0.16	92.75 \pm 0.25
	VGG	87.01 \pm 1.04	90.23 \pm 0.20	91.93 \pm 0.26
	VGG + AS	90.44 \pm 0.52	91.92 \pm 0.06	92.92 \pm 0.08
CIFAR-100	ResNet20	61.05 \pm 0.61	64.53 \pm 0.50	66.97 \pm 0.72
	ResNet20 + AS	63.51 \pm 0.64	65.71 \pm 0.32	66.74 \pm 0.77
	ResNet56	60.51 \pm 1.30	64.99 \pm 0.33	68.16 \pm 0.12
	ResNet56 + AS	64.93 \pm 0.43	67.37 \pm 0.32	69.48 \pm 0.39
	VGG	47.07 \pm 2.35	52.00 \pm 0.68	55.07 \pm 1.05
	VGG + AS	61.73 \pm 0.38	64.18 \pm 0.67	66.07 \pm 1.05

- We present the ensembled accuracies where only partial predictions are retained according to the total uncertainty.
- Adversarial Sampling is still helpful under this scenario.

Dataset	Model	20 % data retained	40 % data retained	60 % data retained	80 % data retained
CIFAR-10	ResNet20	99.82 ± 0.08	99.62 ± 0.14	98.55 ± 0.22	94.45 ± 0.30
	ResNet20 + AS	99.90 ± 0.05	99.74 ± 0.11	99.07 ± 0.15	96.21 ± 0.35
	ResNet56	99.90 ± 0.05	99.75 ± 0.10	98.81 ± 0.23	95.14 ± 0.61
	ResNet56 + AS	99.95 ± 0.00	99.81 ± 0.06	99.21 ± 0.11	96.84 ± 0.48
	VGG	99.88 ± 0.03	99.74 ± 0.09	99.28 ± 0.17	96.54 ± 0.10
	VGG + AS	99.93 ± 0.03	99.79 ± 0.09	99.44 ± 0.06	97.68 ± 0.34
CIFAR-100	ResNet20	96.40 ± 0.50	85.05 ± 1.45	74.09 ± 1.41	64.87 ± 1.26
	ResNet20 + AS	96.68 ± 0.68	87.39 ± 1.59	76.43 ± 1.29	66.53 ± 1.06
	ResNet56	93.88 ± 0.73	80.38 ± 1.29	69.82 ± 2.13	61.09 ± 2.56
	ResNet56 + AS	97.18 ± 0.43	88.09 ± 1.89	77.20 ± 1.85	67.35 ± 1.94
	VGG	86.93 ± 0.28	72.96 ± 0.36	62.96 ± 0.76	54.59 ± 0.82
	VGG + AS	96.32 ± 0.23	85.61 ± 0.59	74.28 ± 0.72	64.79 ± 0.73

- We argue that the randomness of sampling in Bayesian neural networks causes the performance decrease.
- We propose training with Adversarial Distribution as a theoretical solution.
- We further propose Adversarial Sampling as an approximation in practice.
- Extensive experiments validate our proposal.

Thanks
