# Improving Bayesian Neural Networks by Adversarial Sampling

Jiaru Zhang[1]   Yang Hua[2]   Tao Song[1]   Hao Wang[3]   Zhengui Xue[1]   Ruhui Ma[1]   Haibing Guan[1]

[1]Shanghai Jiao Tong University   [2]Queen's University Belfast   [3] Louisiana State University

**AAAI-22**

## Introduction

**Background:** Bayesian neural networks (BNNs) have drawn extensive interest due to the unique probabilistic representation framework. However, Bayesian neural networks have limited publicized deployments because of the relatively poor model performance in real-world applications.

**Goal:** Explore the reason of the relatively poor performance of Bayesian neural networks, and improve the performance by targeted solutions.

**Key Contributions:**
- We argue that the randomness of sampling in Bayesian neural networks causes errors in updating parameters during training and models with poor performance in testing.
- We propose to train Bayesian neural networks with Adversarial Distribution. It can improve the worst performance of the model in multiple samplings and enhance its predictive performance.
- We further propose the Adversarial Sampling method as a practical approximation.
- Verify the theoretical analysis and the effectiveness of the proposed method by experiments under multiple situations.

## Adversarial Sampling

The calculation of $Q_{adv}$ analytically is difficult. We propose an iterative approach, Adversarial Sampling, as an approximation. We first sample each parameter from the original parameter distribution.

$$w_{adv} \sim N(\mu, \sigma^2). \qquad (4)$$

Then we adversarially perturb the parameter $w$ by repeatedly perturb the parameters on the opposite direction of gradient.
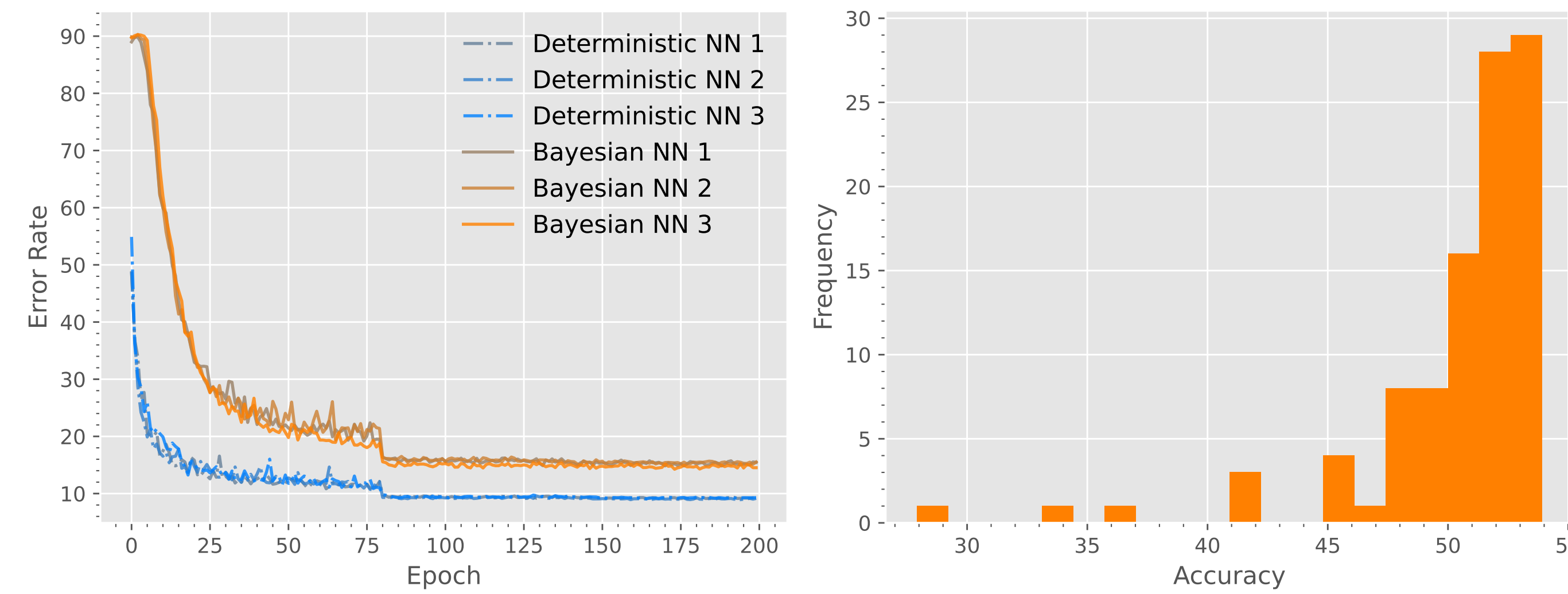
$$w_{adv} = w_{adv} + \alpha \cdot \sigma \cdot \mathrm{sign}\left(\mathrm{grad}\left(w_{adv}\right)\right). \qquad (5)$$

We adjust the scope of the adversarial perturbation using the standard deviation of the parameter $\sigma$, since a parameter with a larger standard deviation has higher randomness in regular sampling.

## 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.



## Training with Adversarial Distribution

Adversarial distribution $Q_{adv}$:

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

Adversarial Loss $\mathcal{L}_{adv}$:

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

Total learning target:

$$\theta = \underset{\theta}{\mathrm{argmin}} \left( (1-\lambda) \cdot \mathcal{L}_p + \lambda \cdot \mathcal{L}_{adv} + \mathcal{L}_r \right), \qquad (3)$$

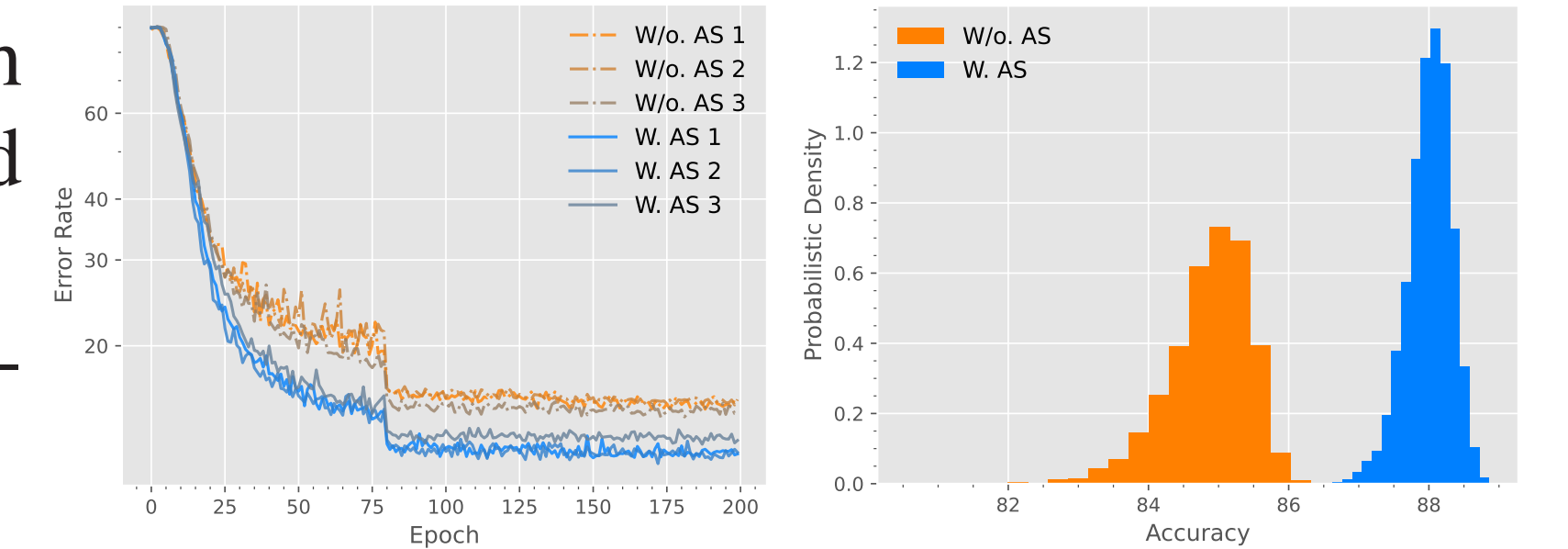- Denoting the iteration times as $N$, the total distance between $w$ and $w_{adv}$ satisfies

$$\|w - w_{adv}\| \leq N \cdot \alpha. \qquad (6)$$

- It satisfies $W[Q_{adv}, Q_\theta] \leq d$ by setting $d = N \cdot \alpha$.
- Many $w_{adv}$s create an approximation of $Q_{adv}$.
- In practice, the parameter $w$ is yielded by a random unit Gaussian noise $\epsilon \sim \mathcal{N}(0,1)$: $w = \mu + \epsilon \cdot \sigma$ with the popularly used reparameterization trick.
- Therefore, we just need to update the random noise $\epsilon$ with the same step size $\alpha$, making Adversarial Sampling simple to implement.

## Experiments & Results

### Verification
- The change trends of models trained with Adversarial Sampling are more stable and steady.
- Models trained without Adversarial Sampling distribute more dispersed.



### Improvement on Model Performance
Models trained with Adversarial Sampling have much higher accuracies.

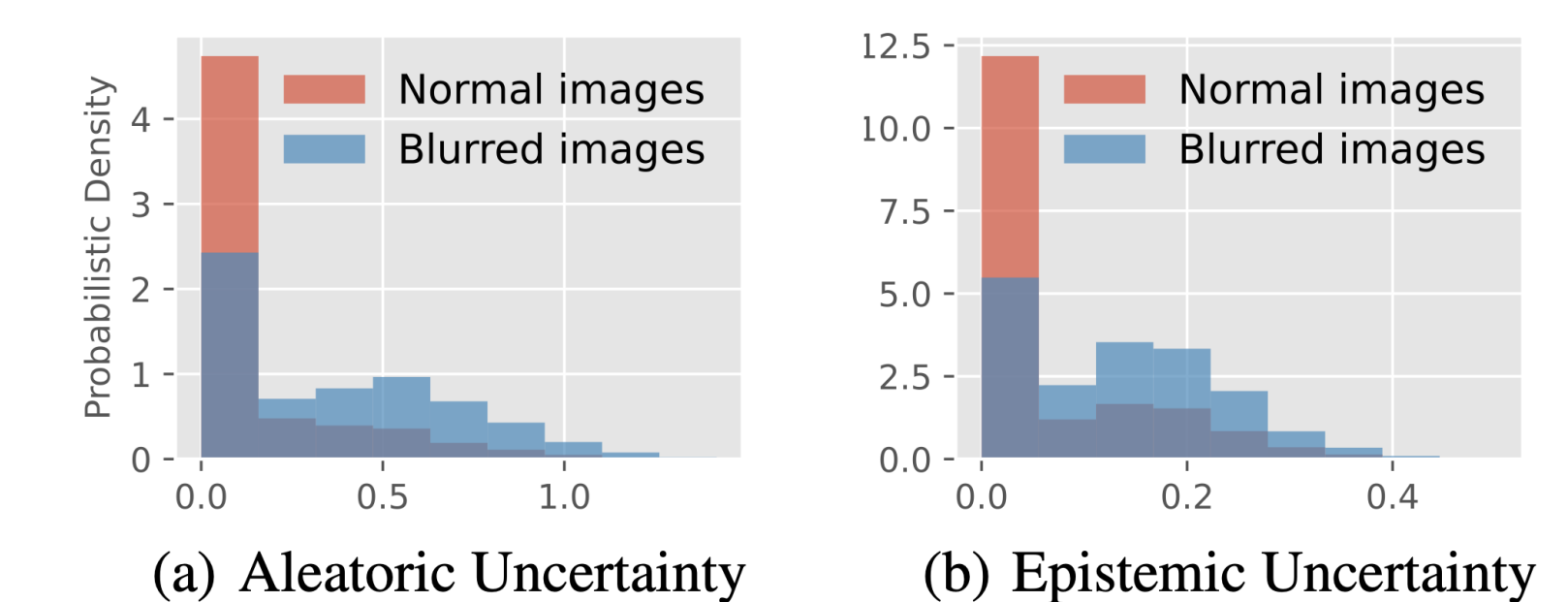| Dataset | Model | Lowest Accuracy | Highest Accuracy | Ensembled Accuracy |
|---|---|---|---|---|
| CIFAR-10 | ResNet20 | 82.73 ± 0.88 | 86.03 ± 0.43 | 87.01 ± 0.65 |
| | ResNet20 + AS | **86.33 ± 0.45** | **88.35 ± 0.51** | **88.76 ± 0.73** |
| | ResNet56 | 82.71 ± 0.55 | 86.84 ± 0.04 | 88.22 ± 0.41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

### Combination with uncertainty estimation
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 | 96.93 ± 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** |

### The Ability of Uncertainty Estimation
Models trained with Adversarial Sampling keep the ability to model uncertainties.



(a) Aleatoric Uncertainty    (b) Epistemic Uncertainty

### Influence of the parameter $\lambda$
Using a suitable $\lambda$ is important.



### Combination with Bayesian Fine-tune
Models trained with the Adversarial Sampling method also perform obviously better compared with original models on this higher baseline.

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

**GitHub Repository:**
AISIGSJTU/AS