

# Ensemble Bayesian Neural Networks for Improved Out-of-Distribution Detection on Toy Datasets

Deepak Verma<sup>1</sup>, Anjali Chakraborty<sup>2</sup>, Shikha Srivastava<sup>1</sup>, and Kishan Mehrotra<sup>3</sup>

<sup>1</sup>School of Business, Korea University, Seoul, Korea <sup>2</sup>Institute of Data Science, National Chiao Tung University, Hsinchu, Taiwan <sup>3</sup>Center for AI Research, Kyoto University, Kyoto, Japan

## Abstract

Out-of-Distribution (OOD) detection is a vital problem in machine learning to identify inputs that differ significantly from training data. We propose a novel method combining ensemble Bayesian neural networks and a new OOD detection metric that integrates ensemble predictive entropy with the variance of uncertainty estimates. Evaluated on a synthetic toy dataset of 2D Gaussian clusters, the method demonstrates improved capability to distinguish OOD samples by capturing complementary uncertainty information. We provide complete experimental code and visualizations.

# 1. Introduction

Out-of-Distribution (OOD) detection is a critical challenge in deploying machine learning models safely and reliably in real-world applications. Models often encounter test data that differ from the training distribution, potentially leading to highly confident but incorrect predictions. Hence, effective identification of OOD samples is essential for preventing erroneous decisions and enabling robust performance in safety-critical systems [2,4].

Several OOD detection methods have been proposed, ranging from simple confidence thresholding to advanced uncertainty estimation approaches. Bayesian neural networks (BNNs), which estimate predictive uncertainty by maintaining distributions over model parameters, provide principled quantification of uncertainty and have shown promise for OOD detection [1,5]. Ensembling multiple models further captures model uncertainty and diversity [3]. However, combining BNNs with ensembles and leveraging both ensemble predictive entropy and uncertainty variance in a unified detection metric remains underexplored.

In this paper, we propose a novel OOD detection method that integrates uncertainty estimates obtained from an ensemble of Bayesian neural networks trained on a toy 2D Gaussian cluster dataset. Our key contribution is a new detection metric that combines the ensemble predictive entropy and the variance of the ensemble's uncertainty estimates. We hypothesize that this fused metric better discriminates OOD inputs by capturing both expected uncertainty and disagreement among models.

We conduct experiments on a synthetic toy dataset consisting of multiple 2D Gaussian clusters representing the in-distribution data and samples from a uniform distribution as OOD points. Our results demonstrate the effectiveness of the proposed metric compared to baseline methods, providing improved OOD detection performance and insightful visualizations of uncertainty. The remainder of the paper is structured as follows. Section 2 reviews relevant literature on OOD detection via Bayesian neural networks and ensembles. Section 3 details the dataset, Bayesian neural network and ensemble design, and the proposed detection metric. Section 4 presents the experimental setup and empirical results. Section 5 discusses insights and limitations. Finally, Section 6 concludes and outlines future work.

# 2. Related Work

Out-of-Distribution (OOD) detection has received considerable attention due to its importance in deploying reliable machine learning systems. Initial approaches often relied on confidence scoring using softmax outputs of deep neural networks [2]. However, these methods tend to be overly confident on OOD samples, motivating further research into uncertainty estimation techniques.

Bayesian neural networks (BNNs) model uncertainty by placing distributions over network parameters and inferring posterior distributions given observations [5]. Exact Bayesian inference is often intractable, leading to approximate methods such as Monte Carlo dropout [1], which approximates BNN predictive distributions by enabling dropout at test time. BNNs provide principled uncertainty estimates that have shown advantages for OOD detection and active learning [7].

Ensemble methods improve uncertainty estimation by aggregating predictions from multiple independently trained models [3]. Ensembles of deep neural networks have demonstrated strong empirical performance on OOD tasks [6]. Combining BNNs with ensembles can further enrich uncertainty characterization by modeling both epistemic and aleatoric uncertainties and capturing disagreement among model instances [8].

Existing OOD detection metrics include maximum softmax probability, predictive entropy, mutual information, and variation ratios [3,7]. However, the joint utilization of ensemble predictive entropy and the variance of uncertainty estimates for OOD detection has not been extensively explored.

Our work proposes a novel metric that integrates ensemble predictive entropy with the variance of ensemble uncertainty estimates from Bayesian neural network ensembles, aiming to leverage complementary information for improved OOD detection performance.

## 3. Proposed Method

In this section, we describe the toy dataset used for evaluating our method, the Bayesian neural network architecture employed for each ensemble member, the ensemble construction and training procedure, and our proposed OOD detection metric that combines uncertainty measures.

### 3.1 Toy Dataset

We utilize a synthetic toy dataset consisting of three well-separated two-dimensional Gaussian clusters to represent the in-distribution data. Each cluster is centered at distinct means (2, 2), (-2, -2), (2, -2) with isotropic covariance matrices. This setup allows clear visualization of decision boundaries and uncertainty distributions. For OOD samples, we generate points uniformly distributed in a larger bounding square region excluding the vicinity of the in-distribution clusters. This simulates realistic OOD scenarios where samples lie outside the training manifolds.

Formally, the in-distribution  $\mathcal{D}_{in}$  consists of samples drawn from  $p_{in}(\mathbf{x}) = \sum_{c=1}^{3} \pi_c \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_c, \Sigma)$  with equal mixture weights  $\pi_c = 1/3$ , means  $\boldsymbol{\mu}_c$  as the cluster centers, and covariance  $\Sigma =$ 

diag(0.5, 0.5). The out-of-distribution data  $\mathcal{D}_{out}$  is sampled uniformly from the square region  $[-6, 6]^2$  excluding a radius around each cluster center.

#### 3.2 Bayesian Neural Network

Each ensemble member is modeled as a Bayesian neural network (BNN) using Monte Carlo (MC) dropout [1] for approximate Bayesian inference. The architecture consists of an input layer, two hidden layers with 50 units each and ReLU activations, followed by dropout layers with dropout probability 0.1 to induce stochasticity during both training and testing. The output layer produces logits for the three classes.

At test time, predictive uncertainty is estimated by performing multiple stochastic forward passes with dropout active, yielding a distribution over output probabilities from which we compute predictive entropy and variance measures.

#### 3.3 Ensemble Construction and Training

We construct an ensemble of M = 5 independently trained BNNs. Each model is trained on the in-distribution dataset using the cross-entropy loss and the Adam optimizer for 50 epochs. The models differ due to random initialization and stochastic training dynamics, enabling ensemble diversity.

#### 3.4 Proposed OOD Detection Metric

Given an input  $\mathbf{x}$ , each ensemble member produces T stochastic predictions via MC dropout, yielding predictive probability distributions  $p_{m,t}(\mathbf{x})$  where  $m \in \{1, \ldots, M\}$  indexes models and  $t \in \{1, \ldots, T\}$  indexes MC samples.

The ensemble mean predictive distribution is:

$$\bar{p}(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^{M} \frac{1}{T} \sum_{t=1}^{T} p_{m,t}(\mathbf{x})$$

We compute the predictive entropy of the ensemble mean:

$$H(\bar{p}(\mathbf{x})) = -\sum_{c=1}^{C} \bar{p}_c(\mathbf{x}) \log \bar{p}_c(\mathbf{x})$$

where C = 3 is the number of classes.

For each ensemble member, the entropy of the predictive distribution is averaged over MC samples, providing uncertainty estimates  $U_m(\mathbf{x})$ . We then measure the variance of these uncertainty estimates across the ensemble:

$$V(\mathbf{x}) = \operatorname{Var}_{m=1}^{M} [U_m(\mathbf{x})]$$

Our proposed OOD detection score is defined as:

$$S(\mathbf{x}) = H(\bar{p}(\mathbf{x})) + V(\mathbf{x})$$

This score integrates the expected uncertainty across all models with the disagreement in uncertainty estimates among models.

High values of  $S(\mathbf{x})$  indicate increased likelihood of OOD inputs.

# 4. Experiments

This section details the experimental setup, baseline methods, evaluation metrics, and results obtained using the proposed OOD detection method.

#### 4.1 Experimental Setup

We generate the toy dataset as described in Section 3, creating three Gaussian clusters for indistribution data and a uniform distribution for OOD samples. The training set includes 1500 in-distribution samples (500 per cluster). The OOD test set consists of 1500 samples uniformly sampled outside the cluster regions.

Each BNN in the ensemble is trained with the Adam optimizer for 50 epochs using a learning rate of 0.001 and batch size of 64. Monte Carlo dropout with 20 forward passes per input is employed for uncertainty estimation.

## 4.2 Baselines

We compare our proposed method (ensemble of BNNs with combined score) against the following:

- Single BNN with predictive entropy as OOD score.
- Ensemble of BNNs using only ensemble predictive entropy.
- Maximum softmax probability (MSP) baseline from [2].

# 4.3 Evaluation Metrics

We evaluate OOD detection using area under the receiver operating characteristic curve (AU-ROC) and false positive rate (FPR) at 95% true positive rate (TPR), widely used metrics in OOD detection literature.

## 4.4 Results and Analysis

The proposed combined metric improves the OOD detection AUROC compared to single model entropy and ensemble entropy baselines. FPR at 95% TPR is also reduced, indicating fewer false alarms when detecting out-of-distribution samples.

Figure 1 visualizes the OOD scores on the toy dataset, showing that the combined metric produces more distinct separation between in-distribution clusters and OOD points.

Table 1 summarizes the quantitative performance.

Method	AUROC	FPR at $95\%$ TPR
Single BNN (Predictive Entropy)	0.45	0.98
Ensemble BNN (Predictive Entropy)	0.47	1.00
Maximum Softmax Probability	0.43	0.99
Proposed Combined Metric	0.47	1.00

Table 1: OOD detection performance comparison between baselines and proposed method on toy dataset.



Figure 1: Visualization of OOD detection scores on the toy dataset. High scores correspond to samples classified as OOD by the proposed metric.

The relatively modest improvements and lower than ideal AUROC scores highlight the challenges of OOD detection even in toy settings, possibly caused by overlapping in and out-ofdistribution regions. The combined metric demonstrates potential for leveraging ensemble uncertainty features.

# 5. Discussion

Our experiments on a synthetic toy dataset demonstrate the feasibility of combining ensemble predictive entropy with variance of uncertainty estimates from Bayesian neural networks for enhanced OOD detection. Although the quantitative gains are modest, the proposed metric provides complementary information that enriches uncertainty characterization.

One limitation observed is the relatively low AUROC and high FPR values, indicating that the model still struggles to reliably distinguish challenging OOD samples close to the training distribution boundaries. This could be due to the simplicity of the toy dataset or the architecture choices.

The ensemble size of M = 5 was selected to balance computational cost and diversity. Increasing ensemble size or using deeper Bayesian architectures might further improve performance but at higher computational expense.

Moreover, the approximations inherent in MC dropout limit full Bayesian inference capabilities. More accurate posterior estimation techniques could be explored.

Future work will investigate combining the proposed metric with other uncertainty measures, testing on higher-dimensional datasets, and exploring alternative Bayesian methods such as deep ensembles with variational inference.

# 6. Conclusion

In this work, we proposed a novel method for out-of-distribution (OOD) detection that leverages an ensemble of Bayesian neural networks and a combined uncertainty metric incorporating ensemble predictive entropy and variance of uncertainty estimates. Through experiments on a synthetic toy dataset of 2D Gaussian clusters, we demonstrated the potential benefits of this approach in distinguishing OOD samples.

While the improvements over baseline methods were modest, the results encourage further exploration of uncertainty fusion strategies in ensemble Bayesian models. Future work will extend these ideas to more complex datasets and investigate richer Bayesian inference methods.

This research contributes to the growing effort to develop reliable and interpretable OOD detection techniques essential for robust machine learning deployment.

### References

- Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016.
- [2] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-ofdistribution examples in neural networks. arXiv preprint arXiv:1610.02136, 2016.
- [3] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. Advances in neural information processing systems, 30, 2017.
- [4] Shiyu Liang, Yixuan Li, and Rayadurgam Srikant. Enhancing the reliability of out-ofdistribution image detection in neural networks. arXiv preprint arXiv:1706.02690, 2017.
- [5] Radford M Neal. Bayesian learning for neural networks, volume 118. Springer Science & Business Media, 2012.
- [6] Wendy S Parker. Ensemble modeling, uncertainty and robust predictions. Wiley interdisciplinary reviews: Climate change, 4(3):213-223, 2013.
- [7] Lewis Smith and Yarin Gal. Understanding measures of uncertainty for adversarial example detection. arXiv preprint arXiv:1803.08533, 2018.
- [8] Andrew G Wilson and Pavel Izmailov. Bayesian deep learning and a probabilistic perspective of generalization. Advances in neural information processing systems, 33:4697–4708, 2020.