
Adaptive Uncertainty Estimation via High-Dimensional Testing on Latent Representations

Supplementary Materials

Anonymous Author(s)

Affiliation

Address

email

1 **Overview.** In the supplementary material, we first present a detailed summary of the datasets in
2 Section 1. Additionally, we provide detailed descriptions of distributions used for BNN training in
3 Section 2. We further present the implementation settings, hyperparameters and settings of baseline
4 methods in Section 3. We also provide the definitions of uncertainty measures used for comparison
5 (Section 4), together with an algorithm of our framework (Algorithm 1).

6 1 Additional Information on Datasets

7 Table 1 presents a summary of the datasets used for experiments. We totally use seven image datasets
to evaluate our method, including nature image and medical image.

Table 1: Summary of Datasets

Dataset	No. Classes	No. Training	No. Testing
MNIST	10	60,000	10,000
Fashion-MNIST	10	60,000	10,000
OMNIGLOT	50	13,180	19,280
SVHN	10	73,257	26,032
CIFAR-10	10	60,000	10,000
CIFAR-100	100	60,000	10,000
DRD	2	50	100

8

9 **Diabetes Retinopathy Detection (DRD).** For this experiment, we define normal samples as healthy
10 (no DR; with label 0), and OOD samples as DR (mild, moderate, severe, or proliferative DR; with
11 labels 1–4). We select 50 healthy images to train the encoder, and compute μ_1 and Σ_1 from these
12 samples using the trained encoder. For testing, we select 50 images as in-distribution data and 50
13 images as the OOD data. All images are resized to 64×64 for computational convenience. We train
14 the encoder with a task to classify whether the input image is left eye or right eye. The examples of
15 healthy and DR are presented in Figure 1.

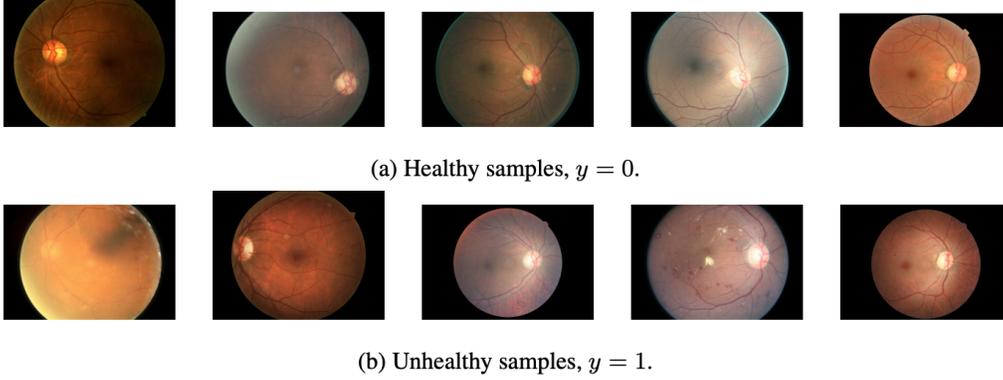


Figure 1: Health and unhealthy (DR) samples from the Diabetes Retinopathy Detection (DRD) dataset [2]

16 2 Multivariate Gaussian Distribution

17 The Multivariate Gaussian is crucial for the approximation of vanilla BNN [5], we provide the
 18 formal definition and its important properties in this section. The density of a multivariate Gaussian
 19 distribution is defined as

$$p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{\frac{n}{2}} |\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\},$$

20 where $\boldsymbol{\mu} \in \mathbb{R}^p$ is a p -dimensional mean vector and $\boldsymbol{\Sigma} \in \mathbb{R}^{p \times p}$ is the covariance matrix.

21 The KL divergence between two multivariate normal distributions $\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ and $\mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$ is
 22 given by

$$\text{KL}(\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \parallel \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)) = \frac{1}{2} \left[\log \frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} - p + \text{tr}\{\boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1\} + (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}_2^{-1} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) \right].$$

23 3 Baseline Methods and Implementation Details

24 **Implementation Details.** We present additional implementation details and hyperparameter settings.
 25 We first provide the key settings and adaptations applied to the baseline methods for reproducibility.
 26 We follow the default settings for other fine-grained parameters (e.g., learning rates).

27 The proposed method is implemented in Python with *Pytorch* library on a server equipped with four
 28 NVIDIA TESLA V100 GPUs. The dropout ratio of each dropout layer is selected as 0.2. All models
 29 are trained with 100 epochs with possible early stopping. We use the *Adam* optimizer to optimize the
 30 model with a learning rate of 5×10^{-5} and a weight decay of 1×10^{-5} . Data augmentations such as
 31 color jittering and random cropping and flipping are applied as a regularization measure.

32 **Hyperparameter Settings.** The hyperparameter settings for BNN training and ARHT testing are

- 33 • Prior of mean of weights — sample from $\mathcal{N}(-3, 0.01)$
- 34 • $s = 5$
- 35 • $n_2 = 300$
- 36 • $\lambda_0 = 0.01$

37 **Additional Settings of Baseline Methods.** We further introduce the experimental settings of baseline
 38 methods:

- 39 • Deep ensembles [7]: Set the number of ensembles as 5.
- 40 • MCDropout [3]: Set the dropout ratio as 0.2 for both training and inference.
- 41 • Kendall and Gal: Set the number of inference weights samples as 20.
- 42 • Detectron [4]: Set the number of runs as 100.
- 43 • PostNet [1]: Because the original codes operate on the features extracted from the standard
- 44 datasets. Their method cannot be generalized to new datasets (e.g., SVHN) as the data
- 45 processing codes are not provided.

46 **4 Uncertainty Measures**

47 We describe the uncertainty measures used in the OOD misclassification task in this section. The
 48 definitions are well-known and summarized by Malinin and Gales [8],

- 49 • Entropy:

$$H[p(\boldsymbol{\mu}|\mathcal{D})] = - \int_{S^{K-1}} p(\boldsymbol{\mu}|\mathcal{D}) \ln p(\boldsymbol{\mu}|\mathcal{D}) d\boldsymbol{\mu},$$

50 where S^{K-1} is the supporting set, and $\boldsymbol{\mu}$ is the predictive class probability which is assumed
 51 to follow Dirichlet distribution.

- 52 • Maximum probability: we take the maximum predicted probability \mathcal{P} from all classes as the
 53 confidence score,

$$\mathcal{P} = \max_c P(w_c|\mathcal{D}).$$

54 where $P(w_c|\mathcal{D})$ is the predictive probability of class c .

- 55 • Differential entropy:

$$I[y, \boldsymbol{\mu}|\mathcal{D}] = - \int_{S^{K-1}} p(\boldsymbol{\mu}|\mathcal{D}) \ln p(\boldsymbol{\mu}|\mathcal{D}) d\boldsymbol{\mu}$$

- 56 • Accuracy: the fraction of correct predictions to the total number of ground truth labels.
- 57 • F-1 score: The F-1 score for each class is defined as

$$\text{F-1 score} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

58 where ‘recall’ is the fraction of correct predictions to the total number of ground truths
 59 in each class and precision is the fraction of correct predictions to the total number of
 60 predictions in each class.

- 61 • AUROC: the area under the receiver operating curve (ROC) which is the plot of the true
 62 positive rate (TPR/Recall) against the false positive rate (FPR).
- 63 • AUPR: the area under the precision-recall curve. Note that the AUPR for binary classification
 64 is sensitive to the distribution of positive and negative classes. Hence, the higher AUPR
 65 does not necessarily imply a better model performance.

66 **5 Algorithm**

67 Algorithm 1 demonstrates the detailed workflow of our proposed uncertainty estimation framework.

Algorithm 1 Our proposed uncertainty estimation framework

Input:
The prior distribution of weights of BNN encoder $\pi(\boldsymbol{\theta}) \sim \mathcal{N}(0, \mathbf{I})$;
Training data $\mathcal{D}_{tr} = \{\mathbf{x}_i, y_i\}_{i=1}^{N_{tr}}$;
Testing data $\mathcal{D}_{te} = \{\mathbf{x}_i, y_i\}_{i=1}^{N_{te}}$;
Hyperparameters μ_0, ρ_0, n_2 ;
Initial variational posterior distribution $q(\boldsymbol{\theta}) \sim \mathcal{N}(\boldsymbol{\mu}, \log(1 + \exp(\boldsymbol{\rho})))$ with initial parameters $\boldsymbol{\mu} = \mu_0 \mathbf{1}$ and $\boldsymbol{\rho} = \rho_0 \mathbf{1}$

Output: The uncertainty scores

- 1: **for** (\mathbf{x}_i, y_i) in \mathcal{D}_{tr} **do** ▷ Train BNN encoder
- 2: Draw weight sample $\boldsymbol{\theta}$ from $q(\boldsymbol{\theta})$
- 3: $\hat{y}_i = f_{\boldsymbol{\theta}}(\mathbf{x}_i)$ ▷ Forward propagation
- 4: Compute task-specific loss \mathcal{L}_{obj}
- 5: Compute the $\text{KL}(q||\pi)$ and hence the ELBO
- 6: Backpropagate the ELBO to update $\boldsymbol{\mu}$ and $\boldsymbol{\rho}$
- 7: **end for**
- 8: Compute $\mu_1 \in \mathbb{R}^p, \boldsymbol{\Sigma}_1 \in \mathbb{R}^{p \times p}$ ▷ Get summary statistics of training
- 9: **for** (\mathbf{x}_i, y_i) in \mathcal{D}_{tr} **do** ▷ OOD Detection
- 10: Compute $\mu_2 \in \mathbb{R}^p, \boldsymbol{\Sigma}_2 \in \mathbb{R}^{p \times p}$
- 11: Compute the pooled sample covariance matrix by Eq. (1)
- 12: Compute ARHT by Eq. (4) as the uncertainty score
- 13: Detect OOD samples using the uncertainty score under famil-wise adjusted threshold
- 14: **end for**

68 **References**

- 69 [1] Bertrand Charpentier, Daniel Zügner, and Stephan Günnemann. Posterior network: Uncertainty
70 estimation without ood samples via density-based pseudo-counts. *Advances in Neural Information*
71 *Processing Systems*, 33:1356–1367, 2020.
- 72 [2] Angelos Filos, Sebastian Farquhar, Aidan N Gomez, Tim GJ Rudner, Zachary Kenton, Lewis
73 Smith, Milad Alizadeh, Arnoud de Kroon, and Yarin Gal. A systematic comparison of bayesian
74 deep learning robustness in diabetic retinopathy tasks. *arXiv preprint arXiv:1912.10481*, 2019.
- 75 [3] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model
76 uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059.
77 PMLR, 2016.
- 78 [4] Tom Ginsberg, Zhongyuan Liang, and Rahul G Krishnan. A learning based hypothesis test for
79 harmful covariate shift. *arXiv preprint arXiv:2212.02742*, 2022.
- 80 [5] Laurent Valentin Jospin. Hands-on bayesian neural networks-a tutorial for deep learning users.
81 2020.
- 82 [6] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for
83 computer vision? *Advances in neural information processing systems*, 30, 2017.
- 84 [7] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predic-
85 tive uncertainty estimation using deep ensembles. *Advances in neural information processing*
86 *systems*, 30, 2017.
- 87 [8] Andrey Malinin and Mark Gales. Predictive uncertainty estimation via prior networks. *Advances*
88 *in neural information processing systems*, 31, 2018.