



Problem: Detecting Out-of-Distribution





[Input]

• Detect whether a test sample is from in-distribution P_{in} (i.e., training distribution) or out-of-distribution P_{out} .

Related work:

- Threshold-based detector [Hendrycks et al., 2016] defined a confidence score as the maximum value of prediction and classifies it as in-distribution if the confidence score is above some threshold.
- ODIN [Liang et al., 2018] further improved the performance using temperature scaling and input pre-processing.

Overconfidence issue: Deep neural networks (DNNs) are typically overconfident in their predictions [Lakshminarayanan et al., 2017]:



Contribution 1: Training Confidence-calibrated Classifier

Training method for detecting out-of-distribution.

• Confidence loss: additionally minimizing the KL divergence.



Experimental results on simple CNNs (2 Conv + 3 FC): some explicit out-ofdistribution samples (denoted by "seen") are given in training time.



Figure 1: (a)/(b) The x-axis and y-axis represent the maximum prediction value and the fraction of images receiving the corresponding score, respectively. (c) Receiver operating characteristic (ROC) curve plots true positive rate (TPR) against false negative rate (FPR).

- In case of confidence loss (Figure 1(b)), there exists a better separation between in- and out-of- distributions.
- Optimizing the cross entropy has a higher FPR than other ones doing the confidence loss to have a same TPR.

Training Confidence-calibrated Classifiers for Detecting Out-of-Distribution Samples

Kimin Lee¹, Honglak Lee^{2,3}, Kibok Lee², Jinwoo Shin¹ ¹Korea Advanced Institute of Science and Technology (KAIST)

If score > ϵ : In-distribution

Contribution 2: GAN for Out-of-Distribution

Issue: the number of out-of-distribution training samples might be almost infinite to cover the entire space. **Intuition:** out-of-distribution samples close to in-distribution are more effective in improving the detection performance.

Experiments on binary classification task:







(a) Out samples from (b) Decision boundary (c) Out samples close to (d) Decision boundary in the case of (c) entire space in the case of (a) in-distribution Figure 2: (a)/(b) Classifier becomes overconfident when out samples are from entire space. (c)/(d) It shows confidence-calibrated predictions if out samples are close to in-distribution.

• Training out-of-distribution samples nearby in-distribution could be more effective in improving detection performances.

New generative adversarial network (GAN): generating most effective samples from out-of-distribution.

 $\min_{G} \max_{D} \quad \mathbb{E}_{P_{G}(\mathbf{x})} \left[KL \left(\mathcal{U} \left(y \right) \parallel P_{\theta} \left(y | \mathbf{x} \right) \right) \right]$ (a) Forcing G to produce low-density samples

(b) Original GAN loss for generating samples close to in-distribution

Comparison of original GAN and proposed GAN:





(a) Samples from the orig- (b) Samples from the proinal GAN

posed GAN

Figure 3: In-distribution samples in (a)/(b) and (c)/(d) are drawn from a mixture of two Gaussian distributions and MNIST dataset, respectively.

- For our method, we use a pre-trained classifier.
- The proposed GAN can produce the samples nearby the low-density boundary of the in-distribution space.

Contribution 3: Joint Training Method

Jointly optimizing confidence-calibrated classifier and GAN as follows:

 $\mathbb{E}_{P_{\text{in}}(\widehat{\mathbf{x}},\widehat{y})} \left[-\log P_{\theta} \left(y = \widehat{y} | \widehat{\mathbf{x}} \right) \right] + \beta \mathbb{E}_{P_{G}(\mathbf{x})} \left[KL \left(\mathcal{U} \left(y \right) \parallel P_{\theta} \left(y | \mathbf{x} \right) \right) \right]$

(c) Standard classification loss $+\mathbb{E}_{P_{\text{in}}(\widehat{\mathbf{x}})}\left[\log D\left(\widehat{\mathbf{x}}\right)\right]+\mathbb{E}_{P_{G}(\mathbf{x})}\left[\log\left(1-D\left(\mathbf{x}\right)\right)\right].$

(e) Original GAN loss

• Confidence loss: (c) + (d) / Proposed GAN loss: (d) + (e)

²University of Michigan ³Google Brain









(c) Images from the original GAN

(d) Images from the proposed GAN

(d) KL divergence term for confident prediction

Experiments: Effects of Joint Confidence Loss

Detection performance of threshold-based detector using VGGNet: • True negative rate (TNR) at 95% true positive rate (TPR). • Area under receiver operating characteristic curve (AUROC): TPR vs FPR. • Detection accuracy: best classification accuracy over all thresholds.



• Our method outperforms all baseline methods in all cases.



Figure 5: Performances of the baseline detector and ODIN detector [Liang et al. 2018]. • Our method can be utilized with ODIN [Liang et al., 2018].

Visual interpretations of trained models:



Figure 6: Guided gradient maps of predicted class with respect to the input image.

