

Detecting Out-of-Distribution Inputs in Deep Neural Networks Using an Early-Layer Output

Vahdat Abdelzad*
vabdelza@uwaterloo.ca

Krzysztof Czarnecki*
kczarnek@gsd.uwaterloo.ca

Rick Salay*
rsalay@gsd.uwaterloo.ca

Taylor Denouden†
tadenoud@uwaterloo.ca

Sachin Vernekar†
sverneka@uwaterloo.ca

Buu Phan*
btphan@uwaterloo.ca

Abstract

Deep neural networks achieve superior performance in challenging tasks such as image classification. However, deep classifiers tend to incorrectly classify out-of-distribution (OOD) inputs, which are inputs that do not belong to the classifier training distribution. Several approaches have been proposed to detect OOD inputs, but the detection task is still an ongoing challenge. In this paper, we propose a new OOD detection approach that can be easily applied to an existing classifier and does not need to have access to OOD samples. The detector is a one-class classifier trained on the output of an early layer of the original classifier fed with its original training set. We apply our approach to several low- and high-dimensional datasets and compare it to the state-of-the-art detection approaches. Our approach achieves substantially better results over multiple metrics.

1 Introduction

Deep Neural Networks (DNNs) are an indispensable part of the current generation of software systems for Autonomous Vehicles (AV), the Internet of Things (IoT), and medical diagnosis. We can observe the power of DNNs specifically in the advanced deep models developed for vision and speech recognition systems (van den Oord et al. 2016; Krizhevsky, Sutskever, and Hinton 2015; He et al. 2015). Classification is one of the tasks that humans perform regularly and deep models sometimes outperform humans in doing this task.

Deep classifiers generalize well when they are given inputs drawn from the same distribution as the training data, referred to as in-distribution (ID). In practice, inputs can be drawn from the in-distribution or other distributions, however. Modern deep networks tend to predict such out-of-distribution (OOD) inputs as an ID class with high confidence (Nguyen, Yosinski, and Clune 2014).

This misbehavior can hinder the adoption of deep classifiers in safety-critical systems. For example, a classifier trained to classify vehicles may misclassify a vulnerable road user on a recumbent bike, if not in the training set,

as a car. Therefore, classifiers need to be enhanced with mechanisms that allow distinguishing ID and OOD inputs. The problem of detecting OOD inputs has been studied extensively in different domains under various names such as outlier and novelty detection (Pimentel et al. 2014; Chandola, Banerjee, and Kumar 2009).

There are several approaches proposed to detect OOD inputs for deep models (Hendrycks and Gimpel 2017; DeVries and Taylor 2018; Liang, Li, and Srikant 2018; Lee et al. 2018b). The baseline approach is max-softmax proposed by Hendrycks and Gimpel (2017), which uses a threshold over the predicted softmax class probability to detect OOD inputs. This approach does not put any assumption over the architecture of deep models or use samples of OOD inputs for detection. It can also be applied to already trained models. However, it does not have satisfactory performance. Other approaches may have better performance than max-softmax, but they constrain the architecture of deep models, put assumptions over the distribution of deep features, cannot be applied to already trained models, or need samples of OOD inputs in order to have better performance (DeVries and Taylor 2018; Liang, Li, and Srikant 2018; Lee et al. 2018b). Therefore, providing a detection approach which is free of those constraints and outperforms max-softmax with a good margin is still a challenge. In this paper, we concentrate specifically on detection approaches that do not require the classifier to be retrained or re-designed, since deep classifiers may be very costly to train, or have specific constraints over architecture.

Although the aforementioned approaches use different techniques to detect OOD inputs, they all essentially rely on features extracted by the penultimate layer of deep classifiers. This layer has been trained to extract features that are important to separate the ID classes of inputs. Figure 1a shows a two-dimensional representation of the features extracted from the penultimate layer of a ResNet model trained on CIFAR-10 (Krizhevsky 2009). The features belong to the test datasets of CIFAR-10 and TinyImageNet (Deng et al. 2009), which are considered ID and OOD, respectively. There are ten small clusters (i.e., manifolds) of features that represent different classes (i.e., ID inputs) of CIFAR-10. There is also a large cluster of features that rep-

*Department of Electrical and Computer Engineering, University of Waterloo

†Department of Computer Science, University of Waterloo

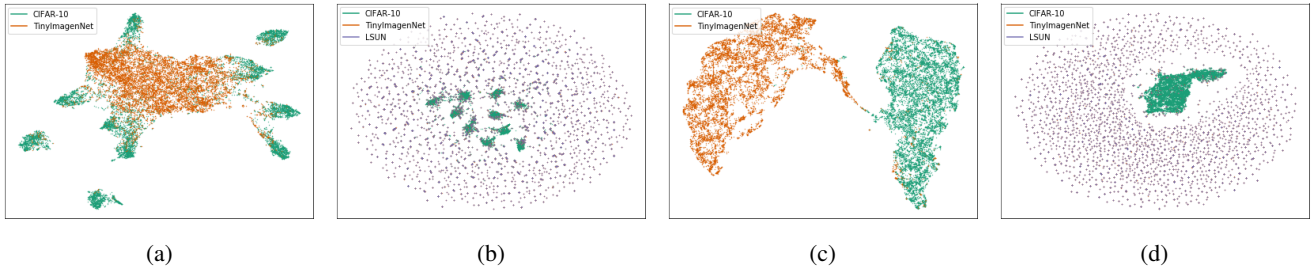


Figure 1: Two-dimensional representations of features extracted from a ResNet model trained on CIFAR-10. a) Features extracted from the penultimate layer for the test set of CIFAR-10 and TinyImageNet datasets. b) Similar to (a) but LSUN was added as the second OOD dataset. c) Similar to (a) but features have been extracted from a well-chosen layer. d) Similar to (b) but features have been extracted from a well-chosen layer. Visualizations use UMAP (McInnes et al. 2018).

resents OOD inputs. In this case, the detection problem manifests itself as separating the distribution of OOD features from the distributions of the different classes (i.e., ten different distributions). This can be challenging even when a distance function or a learning method is used to separate them. Figure 1b depicts the features extracted for CIFAR-10, TinyImageNet, and LSUN datasets (Yu et al. 2015). Here, it becomes much harder to separate ID and OOD inputs. Hein, Andriushchenko, and Bitterwolf (2018) also demonstrate that softmax-based approaches cannot avoid assigning a high-class probability for inputs that are far from the training distribution (i.e., OOD inputs). It is deemed as an inherent problem of DNNs.

We propose a new approach that is not based on class probability (the output of the softmax layer), it significantly outperforms the-state-of-the-art approaches, does not need OOD samples for training the detection model, and does not require retraining the classifier. Our main insight is that there is a latent space in which the distributions of ID and OOD datasets are well separated (regardless of the distribution of classes) and the transformation function to such a latent space has already been approximated well by one of the classifier’s layers. This insight allows learning a detector based on features in this space to separate ID and OOD inputs.

Figure 1c shows a two-dimensional representation of features in such a latent space for the same model and datasets used in Figure 1a. There are two different clusters, each one representing features of ID and OOD datasets, respectively. Being able to represent the entire ID distribution as a manifold in a space can result in better separation. Figure 1d shows features in such a space for an ID and two OOD datasets (similar to Figure 1b). Here, the ID features are again bounded well whereas OOD features are distributed around. In our approach, we find such a latent space and learn a secondary model (which is trained in a few minutes) to separate distributions of features for ID and OOD datasets. We apply our approach to several low- and high-dimensional datasets and then compare it to the baseline and state-of-the-art approaches. The results demonstrate significant improvement for multiple detection metrics.

The rest of paper is organized as follow. In Section 2, we briefly go over related work. Then, we introduce our ap-

proach in Section 3, including how to detect OOD inputs and how to find suitable latent space. In Section 4, we present our experimental results. Finally, we conclude and suggest directions for future work in Section 5.

2 Related work

Hendrycks and Gimpel (2017) propose a baseline approach to detect OOD inputs, called max-softmax, and a set of metrics to evaluate OOD detectors. The approach relies on the observation that ID inputs tend to have a higher predicted softmax class probability compared to OOD inputs. Therefore, a threshold over the predicted softmax class probability should allow separating ID and OOD inputs. Despite being a simple and easy-to-use approach to detect OOD inputs, the approach does not have satisfactory performance, especially for critical applications.

Liang, Li, and Srikant (2018) propose an approach called ODIN that improves on max-softmax by incorporating two extra components: temperature scaling and input preprocessing. Temperature scaling (Guo et al. 2017) is used to calibrate the softmax output of a network, and input preprocessing helps to increase its maximum output for ID inputs. Although preprocessing improves the detection of ID inputs, it needs *access to OOD samples in advance* to fine-tune the perturbation magnitude used in preprocessing during inference. In practice, it is possible to have access to some OOD samples, but it is hard or impossible to have access to samples from all OOD datasets. Input preprocessing also requires extra processing time (two forward and one backward passes over the model), which can be an issue for real-time systems.

The performance of the softmax-based approaches (such as ODIN) depends highly on how the predicted softmax class probability varies for ID and OOD inputs. A higher class probability for an input indicates a higher chance of considering it as ID. Lee et al. (2018a) further improve upon ODIN and propose an approach that forces deep classifiers to output close to uniform distribution for OOD samples. The main idea is to jointly train a classifier and a generator: the generator is trained to produce samples at the boundary of the data manifold, which are considered OOD samples; the classifier is trained using a specially designed

loss function that encourages the classifier to assign uniform class probabilities to these generated OOD samples. This approach requires expensive retraining for an existing classifier and relies on the assumption that generated samples cover the entire boundary of the data manifold, which is difficult to achieve for high-dimensional data (Vernekar et al. 2019).

DeVries and Taylor (2018) utilize uncertainty to detect OOD inputs. They assume that a classifier is more confident about its prediction when the input is ID. Therefore, they proposed to retrain a classifier to output also confidence estimates for each input. Then, the confidence score is used to differentiate between ID and OOD inputs. Others (Shalev, Adi, and Keshet 2018; Vyas et al. 2018) suggested to use ensembles to calculate confidence estimates. Similar to the approach proposed by Lee et al. (2018a), these approaches cannot be applied to already trained models.

MC-dropout (Gal and Ghahramani 2016) is another technique to measure uncertainties of models. This technique relies on multiple inferences to calculate uncertainties and may not be practical in real-time systems. Furthermore, it can only be applied to models that have been trained with dropout layers. Geifman and El-Yaniv (2017) used MC-dropout to draw the risk-coverage curve, which is similar to the concept of OOD detection, for already trained models. They reported that MC-dropout performance is similar to max-softmax. Moreover, uncertainty-based approaches are well-suited for detecting confusing inputs existing near to the boundaries of classes that a model trained for. This affects negatively the effectiveness of these approaches for inputs far from the training distribution.

Using generative models is also another way to detect OOD inputs (Denouden et al. 2018; Pidhorskyi et al. 2018). These approaches usually rely on either the reconstruction error or an estimation of density using the latent representations or a combination of the two. However, these approaches are found to give a higher likelihood to some of the OOD datasets (Nalisnick et al. 2019). Furthermore, generative models tend to scale poorly with the dimension of the dataset. For example, a reconstruction-based OOD detection approach performs well for MNIST, but it becomes difficult to train a generative model to reliably capture the latent manifold of high dimensional datasets (Wang et al. 2017).

Lee et al. (2018b) proposed an approach in which they obtain the class conditional Gaussian distributions with respect to features of the deep models under Gaussian discriminant analysis. This results in a confidence score based on Mahalanobis distance (MD) that is used to detect OOD inputs. They propose two approaches (without considering input preprocessing): a) MD over features before the logits layer; b) an ensemble model based on MDs for all layers. The former performance is similar to ODIN as indicated in their original paper (Table 1 in Lee et al. (2018b)). The latter trains a regression model with OOD samples for each OOD dataset. In fact, for n OOD datasets, n regression models are trained. This is a significant limitation of their approach because the regression models are biased toward the particular OOD distributions on which they are trained. Lee et al. (2018b) also used adversarial samples to train the regression

models. This resulted in a reduction in the detection performance (shown by the variation in the performance on the right-hand side of Table 2 in Lee et al. (2018b)), which is a sign of biases. Last, for both approaches (a & b) proposed by Lee et al. (2018b) distributions of features for each class must follow the multivariate Gaussian distribution. There is no guarantee that distributions of features will satisfy this assumption.

Our approach detailed in Section 3, although seemingly similar to the approach by Lee et al. (2018b) because of using deep features, differs in several important ways: i) it does not use features in the logits layer. We demonstrated in Figures 1c and 1d that features in such a layer are not appropriate to separate ID and OOD; ii) it does not need to have access to OOD samples. *We always train one model (without OOD or adversarial samples) and this single model applies to all OOD datasets*; iii) it finds the appropriate feature space in which ID and OOD inputs are well-separated and thus has no need to rely on ensembles; iv) it does not force the feature distributions to follow the multivariate Gaussian distribution.

3 Proposed solution

In this section, we describe different elements of our approach, including how detection is performed for a specific latent space, how such a latent space is found, and how the approach can be enhanced using input preprocessing.

OOD detection

Let $Q : \mathbb{R}^n \rightarrow [0, 1]^c$ be a function representing a deep network (i.e., classifier), where $x \in \mathbb{R}^n$ is the input and c is the number of classes. Denote as Q_i the output of network Q for class i , and as $X = \{x_1, \dots, x_m\}$ the training set. Network Q has L layers and the output of layer l (i.e., activation values) is represented by q^l ($q^0 = x$). Indeed, q^l is the representation of input x in a latent space obtained by non-linear transformations from layers 1 to l . Each layer allows extracting unique features related to input x . For example, the final layer of a network extracts features that are important to separate the class of input.

The hypothesis underlying our approach is that there is also a latent space, represented by the output of a layer called *optimal OOD discernment layer (OODL)*, in which represented features are discriminative enough to allow separating distributions of ID and OOD datasets. In particular, we deem that features used to separate the class of input might not be appropriate to decide whether or not an input is OOD. If we learn the probability distribution function of features obtained by the OODL, then it should be feasible to separate ID and OOD inputs. Our experiments presented later confirm this hypothesis for the studied datasets and networks.

Now, let l_o be the OODL and Q^{l_o} be the output of layer l_o for every $x_i \in X$. Then, we could train a classifier named S^{l_o} based on Q^{l_o} features to separate ID and OOD inputs. However, to train a two-class classifier S^{l_o} , access is needed to both ID and OOD features. Although there might be some OOD samples available during the training (e.g., related to a specific OOD distribution), the classifier needs to have samples from all OOD distributions in order to be trained well.

Algorithm 1 Finding the OODL for detecting OOD inputs

Require: t_ds : training dataset set, id_ds : ID dataset, ood_ds : OOD dataset, $\{nu, k\}$: training error and kernel for OSVM.

- 1: Initialize the detection error vector: $E = 0$
- 2: **for** each layer $l \in 1, \dots, L$ **do**
- 3: Extract features of layer l for training dataset:
 $Q^l = \text{extract_q}(t_ds, l)$
- 4: Train a one-class SVM classifier:
 $S^l = \text{oneClassSVM,fit}(Q^l, nu, k)$
- 5: Extract features of layer l for ID data:
 $ID^l = _q(id_ds, l)$
- 6: Extract features of layer l for OOD data:
 $OOD^l = \text{extract_q}(ood_ds, l)$
- 7: Calculate the detection error for layer l :
 $E[l] = \text{cal_det_error}(S^l, ID^l, OOD^l)$
- 8: **end for**
- 9: **return** $\text{argmin}(E)$

Furthermore, having access to specific OOD samples during training will cause classifier S^{l_o} to be biased toward detecting those OOD inputs.

Therefore, we formulate this as a one-class classification (OCC) problem. In the OCC context, most of the training samples are ID, and the inputs at inference time are expected to include both ID or OOD samples. Thus, there is no need to have OOD samples during training. OCC has been studied extensively (Khan and Madden 2014; Perera and Patel 2018) and most existing classification methods can be used for this purpose. In this paper, we use One-class Support Vector Machine (OSVM) (Tax and Duin 1999), which is a commonly used one-class classification algorithm. OSVM outputs a score for a given input x , and thresholding over that score allows us to detect OOD inputs. Our detection mechanism is defined as follows.

$$O_{l_o}(x; \delta) = \begin{cases} 0 & \delta \geq S^{l_o}(q^{l_o}(x)) \\ 1 & \text{otherwise} \end{cases} \quad (1)$$

O_{l_o} is the detection function based on features in the OODL l_o and δ is the detection threshold. When the output of classifier S^{l_o} for the features of input x in layer l_o is greater than δ , input x is ID, otherwise, it is OOD. When layer l_o is a fully-connected layer, q^{l_o} is the exact output of layer l_o . However, when layer l_o is a convolutional layer and q^{l_o} is the exact output of layer l_o , q^{l_o} becomes high-dimensional. This can have negative effect on the performance of OSVM, because OSVM performs better for low-dimensional data. Therefore, we calculate the mean of each channel to reduce the dimension of q^{l_o} . Precisely, let $f^{l_o} \in \mathbb{R}^{w \times h \times d}$ be the feature map of convolutional layer l_o , where w , h , and d are width, height, and depth, respectively. Let $f_{ijk}^{l_o}$ be the (i, j, k) -th element of f^{l_o} , then, $q^{l_o} = (q_k^{l_o}) \in \mathbb{R}^d$ is given by

$$q_k^{l_o} = \frac{1}{w \times h} \sum_{i=1}^w \sum_{j=1}^h |f_{ijk}^{l_o}|. \quad (2)$$

Note that traditional OCC methods such as OCSVM perform poorly on high-dimensional data. In our approach, we propose OODL, which allows such methods to perform well on high-dimensional data. An additional benefit is that we can leverage an existing model trained for the original classification task, which is more efficient than training an OCC model from scratch using the ID samples. To the best of our knowledge, no one has proposed using such a setup for OOD detection before.

Finding optimal OOD discernment layer (OODL)

Our approach requires finding the OODL l_o for a given network Q . To do so, we use an OOD dataset to measure the detection error (defined in Section 4) for layers of network Q based on Equation (1). The layer with the minimum detection error is then selected as the OODL.

Algorithm 1 shows how the selection is performed. t_ds is the training dataset of network Q that is used to extract features for layer l . S^l is then trained on the features of layer l according to parameters nu and k that are training error and kernel type, respectively. id_ds and ood_ds are ID and OOD datasets used to measure the detection error based on S^l . id_ds is chosen to be the test set of network Q , and ood_ds can be an arbitrary OOD dataset.

Using an OOD dataset to find the OODL might indicate that such a layer would be biased toward features allowing a better separation of OOD inputs coming from the selected OOD dataset. However, we experimentally found that the choice of an OOD dataset does not affect the approximate position of the OODL. For example, Figure 2 shows the detection error measured for residual layers of a ResNet model trained on CIFAR-10 for different OOD datasets (TinyImageNet, LSUN, iSUN(Xu et al. 2015), and SVHN (Netzer et al. 2011)). As seen, the OODL stays the same ($l_o = 13$) regardless of which OOD dataset is chosen to measure the detection error. *Also note that the OOD dataset is not used during training S^{l_o} .* It is used only to calculate the detection error to select the OODL.

However, the OODL varies based on the ID dataset and deep model. For example, a ResNet model trained on CIFAR-10 has a different OODL than one trained on CIFAR-100 even though the architecture is fixed. Furthermore, the OODL is always one of the low-level layers. This may be associated with the fact that high-level layers are customized to extract features that are useful to separate the classes of input and not to separate the distribution of ID dataset. In other words, high-level features cover unique cases (related to classes), whereas the low-level features cover general cases (related to the distribution) (Zintgraf et al. 2017). Figure 1 demonstrates this perspective using two-dimensional representations of extracted features for the penultimate and OODL of the ResNet model. As seen, the features extracted by the penultimate layer are helpful to separate the classes of input, whereas features of the OODL are useful to separate the distribution of the ID dataset.

Input preprocessing

As indicated in Section 1, input preprocessing is a technique used by ODIN and related approaches to increase the pre-

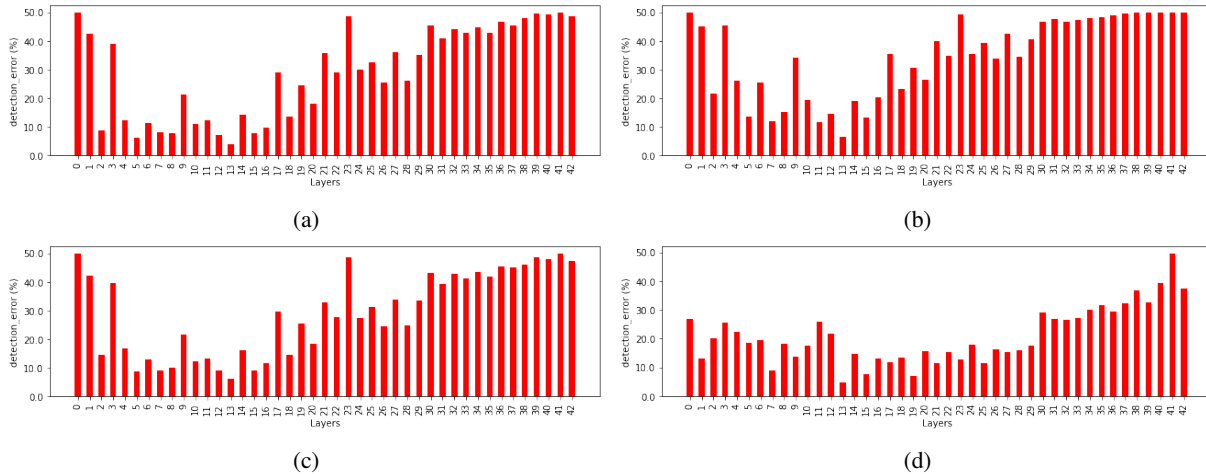


Figure 2: The detection error for residual layers of a ResNet model trained over CIFAR-10. a) The detection error for TinyImageNet. b) The detection error for LSUN. c) The detection error for iSUN. d) The detection error for SVHN. All experiments use the test set of CIFAR-10 as ID dataset.

dicted softmax class probability and thus improve the detection of ID inputs. In our approach, we do not use the output of softmax, but input preprocessing can increase feature values for ID inputs, so that the OSVM classifier can better detect ID inputs. Therefore, we can additionally exploit input preprocessing with our approach; however, Section 4 shows that our approach outperforms other approaches even without input preprocessing.

To implement input preprocessing, we first calculate the pre-processed sample x' , for each input x at the test time, by adding a small perturbation and then use the features of x' for detection. Input x' is obtained as follows

$$x' = x - \varepsilon \cdot \text{sign}\left(-\nabla_x \log(\max_i Q_i(x))\right) \quad (3)$$

where ε is the perturbation magnitude, and the perturbation is calculated by back-propagating the gradients of the predicted class probability with respect to input x .

4 Experiments

We apply our approach to several ID and OOD datasets under different learning models. We compare our approach with max-softmax, ODIN, and MD (over the logits layer with input preprocessing) approaches. As already discussed, we focus on approaches that do not require training on OOD samples or retraining the classifier. Using OOD samples gives access to extra knowledge and retraining a classifier may be very costly. However, to respect the established practice in the literature, we include input preprocessing (with access to OOD) as an optional step in our comparison. Due to the page limitation, other comparisons including uncertainty-based approaches are in supplemental material. Our implementation is available online for reproducibility.

ID datasets and models

We evaluate our approach over several ID datasets, including MNIST, CIFAR-10, and CIFAR-100 (Krizhevsky 2009).

MNIST is a dataset of handwritten digits and has 60,000 images in the training set and 10,000 images on the test set. It includes 28×28 grayscale images. We trained a custom convolutional neural network (CNN) for the MNIST dataset with two convolutional layers and two fully connected layers. The model has an accuracy of 99.22% in the test set. The CIFAR-10 dataset has 50,000 and 10,000 images for training and testing, respectively. It includes 32×32 colored images. We trained two models based on VGG-16 by Simonyan and Zisserman (2014) and ResNet (i.e., ResNet44 v1) by He et al. (2015) for CIFAR-10, achieving an accuracy of 93.56% and 92.01%, respectively. CIFAR-100 is similar to CIFAR-10, but it has 100 classes. We also trained VGG-16 and ResNet models for CIFAR-100, with an accuracy of 70.48% and 69.17%, respectively. Test sets of all ID datasets are used to compute metrics defined in Section 4.

OOD datasets

We consider several OOD datasets for our evaluation, including synthetic ones, and use their test sets for computing the metrics. Moreover, we always keep the size of ID and OOD inputs the same during evaluation (by randomly selecting data from the larger dataset to match the number of instances in the smaller dataset). The following are the OOD datasets used for our experiments.

- **Fashion-MNIST (F-MNIST)** is similar to the MNIST dataset, but it includes Zalando’s article images (Xiao, Rasul, and Vollgraf 2017).
- **Omniglot** contains different handwritten characters from 50 different alphabets. The images have been downsampled to 28×28 images (Lake, Salakhutdinov, and Tenenbaum 2015).
- **TinyImageNet** consists of a subset of ImageNet images (Deng et al. 2009) and covers 200 different classes. We downsampled images to 32×32 .

Table 1: Comparison of our approach with the max-softmax, ODIN, and MD approaches for the MNIST datasets. Values are in percentages and \downarrow indicates that lower values are better, while \uparrow indicates that higher values are better. **Bold** indicates the best score.

ID model	OOD	FPR at 95% TPR \downarrow	Detection error \downarrow	AUROC \uparrow	AUPR Out \uparrow	AUPR In \uparrow
max-softmax / ODIN / MD / ours						
MNIST	F-MNIST	6.77/5.75/34.51/ 0.42	5.82/5.32/19.74/ 2.67	98.06/98.77/86.86/ 99.08	98.44/98.85/90.03/ 98.84	97.55/98.71/82.09/ 99.29
CUSTOM-CNN	Omniplot	6.92/4.79/10.15/ 0.0	5.82/4.53/7.55/ 0.0	97.64/98.8/97.35/ 100.0	98.3/99.04/97.91/ 100.0	96.51/98.5/96.71/ 100.0
	Gaussian	1.74/0.13/26.72/ 0.0	1.22/0.15/14.76/ 0.0	98.91/99.96/77.67/ 100.0	99.38/99.97/87.46/ 100.0	96.86/99.91/62.24/ 100.0
	Uniform	3.76/0.8/30.68/ 0.0	2.87/1.04/17.5/ 0.0	98.11/99.74/78.24/ 100.0	98.84/99.81/87.0/ 100.0	95.74/99.58/63.29/ 100.0

- **LSUN** includes 32×32 downsampled images from the Large-scale Scene UNderstanding dataset (Yu et al. 2015).
- **iSUN** includes 32×32 downsampled images of iSUN images (Xu et al. 2015).
- **SVHN** includes real-world images similar to the MNIST dataset (Netzer et al. 2011).
- **Gaussian noise** includes random normal noise with $\mu = 0.5$ and $\sigma = 1$, clipped to $[0, 1]$.
- **Uniform noise** includes random uniform noise between $[0, 1]$.

Evaluation metrics

There are different metrics to evaluate the performance of OOD detection approaches. We adopted the following metrics (Hendrycks and Gimpel 2017; Liang, Li, and Srikant 2018).

- FPR at 95% TPR is the probability of an out-of-distribution (i.e., negative) input being misclassified as in-distribution (i.e., positive) input when the true positive rate (TPR) is as high as 95%. True positive rate is calculated by $TPR = TP/(TP + FN)$, where TP and FN denote true positives and false negatives, respectively. The false positive rate (FPR) is computed by $FPR = FP/(FP + TN)$, where FP and TN denote false positives and true negatives, respectively.
- Detection error calculates the misclassification probability when TPR is 95%. It is equal to $0.5 * (1 - TPR) + 0.5 * FPR$, where we assume that both positive and negative examples have an equal probability of appearing in the evaluation test.
- AUROC is the Area Under the Receiver Operating Characteristic curve. It is interpreted as the probability that a positive example is assigned a higher detection score than a negative example. An ideal OOD detector expects an AUROC score of 100%.
- AUPR is the Area under the Precision-Recall curve. It is a graph reflecting precision equal to $TP/(TP + FP)$ and recall equal to $TP/(TP + FN)$ against each other. The metric AUPR-In and AUPR-Out represent the area under the precision-recall curve where in-distribution or out-of-distribution images are specified as positives, respectively.

OSVM and hyper parameters

To train OSVM classifiers for detection and finding the OODL, we used the *rbf* kernel and training error $nu = 0.001$. The *rbf* kernel gave us the best results in comparison to other kernels such as *linear* or *poly*. We used temperature scale $T = 1000$ for the ODIN approach since it was deemed as the optimal value based in the original paper (Liang, Li, and Srikant 2018). The perturbation magnitude ϵ for our approach and ODIN was optimized to minimize FPR at 95% TPR by having access to randomly selected 20% of OOD datasets. We used F-MNIST and TinyImageNet datasets to find the OODLs. The OODL was fixed for every ID dataset and its associated model. The OODL for the MNIST dataset was the second convolutional layer. The OODLs for VGG-16 trained on CIFAR-10 and CIFAR-100 were the second convolutional layer and the first max-polling layer, respectively. The OODL for ResNet trained on CIFAR-10 and CIFAR-100 were the thirteenth and ninth residual layers, respectively. Details of parameters for training models are in supplemental material.

Detection results

Table 1 shows the comparison of our approach with others for MNIST (a low-dimensional dataset). Our approach gives better results across all the metrics defined in Section 4. The comparison of our approach *with* input preprocessing with CIFAR-10 and CIDAR-100 are listed in Table 2.a. Our approach also gives better results for these two datasets. Furthermore, our approach is capable of fully detecting noise for both low and high dimensional data under different models. Table 2.b compares our approach *without* input preprocessing with others (*with* input preprocessing). As seen, it outperforms other approaches except in the cases of TinyImageNet for CIFAR-100. Our approach still has better AUROC, but the detection error and FPR at 95% TPR are slightly larger than ODIN’s. Interestingly, the MD approach is worse than max-softmax in some cases. Such a result has also been reported by Ren et al. (2019) in their Table 3 results.

5 Conclusion

Detecting OOD inputs of DNNs is an important concern for the application of DNNs in safety-related domains, such as autonomous driving and medical diagnostics. Most of the current OOD detection approaches rely on features extracted by the penultimate layers of DNNs. Such a layer is trained to extract features that are relevant to separate the classes of

- Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; and Fei-Fei, L. 2009. Imagenet: A large-scale hierarchical image database. In *CVPR*, 248–255.
- Denouden, T.; Salay, R.; Czarnecki, K.; Abdelzad, V.; Phan, B.; and Vernekar, S. 2018. Improving reconstruction autoencoder out-of-distribution detection with mahalanobis distance. *CoRR* abs/1812.02765.
- DeVries, T., and Taylor, G. W. 2018. Learning confidence for out-of-distribution detection in neural networks. *arXiv preprint arXiv:1802.04865*.
- Gal, Y., and Ghahramani, Z. 2016. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *Proceedings of The 33rd International Conference on Machine Learning*, 1050–1059.
- Gal, Y.; Islam, R.; and Ghahramani, Z. 2017. Deep bayesian active learning with image data. *CoRR* abs/1703.02910.
- Geifman, Y., and El-Yaniv, R. 2017. Selective classification for deep neural networks. In *NeurIPS*, 4885–4894.
- Guo, C.; Pleiss, G.; Sun, Y.; and Weinberger, K. Q. 2017. On calibration of modern neural networks. In *ICML*.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2015. Deep residual learning for image recognition. *CoRR* abs/1512.03385.
- Hein, M.; Andriushchenko, M.; and Bitterwolf, J. 2018. Why relu networks yield high-confidence predictions far away from the training data and how to mitigate the problem. *CoRR* abs/1812.05720.
- Hendrycks, D., and Gimpel, K. 2017. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *ICLR*.
- Khan, S. S., and Madden, M. G. 2014. One-class classification: taxonomy of study and review of techniques. *The Knowledge Engineering Review* 29(3):345374.
- Krizhevsky, A.; Sutskever, I.; and Hinton, G. E. 2015. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems* 10971105.
- Krizhevsky, A. 2009. Learning multiple layers of features from tiny images.
- Lake, B. M.; Salakhutdinov, R.; and Tenenbaum, J. B. 2015. Human-level concept learning through probabilistic program induction. *Science* 350(6266):1332–1338.
- Lee, K.; Lee, H.; Lee, K.; and Shin, J. 2018a. Training confidence-calibrated classifiers for detecting out-of-distribution samples. In *ICLR*.
- Lee, K.; Lee, K.; Lee, H.; and Shin, J. 2018b. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In *NeurIPS, NIPS'18*, 7167–7177.
- Liang, S.; Li, Y.; and Srikant, R. 2018. Principled detection of out-of-distribution examples in neural networks. In *ICLR*.
- McInnes, L.; Healy, J.; Saul, N.; and Grossberger, L. 2018. Umap: Uniform manifold approximation and projection. *The Journal of Open Source Software* 3(29):861.
- Nalisnick, E.; Matsukawa, A.; Teh, Y. W.; Gorur, D.; and Lakshminarayanan, B. 2019. Do deep generative models know what they don't know? In *ICLR*.
- Netzer, Y.; Wang, T.; Coates, A.; Bissacco, A.; Wu, B.; and Ng, A. Y. 2011. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning*.
- Nguyen, A. M.; Yosinski, J.; and Clune, J. 2014. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. *CoRR* abs/1412.1897.
- Perera, P., and Patel, V. M. 2018. Learning deep features for one-class classification. *CoRR* abs/1801.05365.
- Pidhorskyi, S.; Almohsen, R.; Adjero, D. A.; and Doretto, G. 2018. Generative probabilistic novelty detection with adversarial autoencoders. *CoRR* abs/1807.02588.
- Pimentel, M. A.; Clifton, D. A.; Clifton, L.; and Tarassenko, L. 2014. Review: A review of novelty detection. *Signal Process.* 99:215–249.
- Ren, J.; Liu, P. J.; Fertig, E.; Snoek, J.; Poplin, R.; DePristo, M. A.; Dillon, J. V.; and Lakshminarayanan, B. 2019. Likelihood ratios for out-of-distribution detection. *CoRR* abs/1906.02845.
- Scheffer, T.; Decomain, C.; and Wrobel, S. 2001. Active hidden markov models for information extraction. In *Advances in Intelligent Data Analysis*, 309–318. Berlin, Heidelberg: Springer Berlin Heidelberg.
- Shalev, G.; Adi, Y.; and Keshet, J. 2018. Out-of-distribution detection using multiple semantic label representations. *Advances in Neural Information Processing Systems 31*:7375–7385.
- Shannon, C. E. 1948. A mathematical theory of communication. *The Bell System Technical Journal* 27(3):379–423.
- Simonyan, K., and Zisserman, A. 2014. Very deep convolutional networks for large-scale image recognition. *CoRR* abs/1409.1556.
- Tax, D. M., and Duin, R. P. 1999. Support vector domain description. *Pattern Recognition Letters* 20(11):1191–1199.
- van den Oord, A.; Dieleman, S.; Zen, H.; Simonyan, K.; Vinyals, O.; Graves, A.; Kalchbrenner, N.; Senior, A.; and Kavukcuoglu, K. 2016. Wavenet: A generative model for raw audio. *arXiv:1609.03499* 1–15.
- Vernekar, S.; Gaurav, A.; Denouden, T.; Phan, B.; Abdelzad, V.; Salay, R.; and Czarnecki, K. 2019. Analysis of confident-classifiers for out-of-distribution detection. In *ICLR workshop on SafeML*.
- Vyas, A.; Jammalamadaka, N.; Zhu, X.; Das, D.; Kaul, B.; and Willke, T. L. 2018. Out-of-distribution detection using an ensemble of self supervised leave-out classifiers. In *ECCV*.
- Wang, W.; Wang, A.; Tamar, A.; Chen, X.; and Abbeel, P. 2017. Safer classification by synthesis. *CoRR* abs/1711.08534.
- Xiao, H.; Rasul, K.; and Vollgraf, R. 2017. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms.
- Xu, P.; Ehinger, K. A.; Zhang, Y.; Finkelstein, A.; Kulkar, S. R.; and Xiao, J. 2015. Turkergaze: Crowdsourc-

ing saliency with webcam based eye tracking. *CoRR* abs/1504.06755.

Yu, F.; Zhang, Y.; Song, S.; Seff, A.; and Xiao, J. 2015. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *CoRR* abs/1506.03365.

Zintgraf, L. M.; Cohen, T. S.; Adel, T.; and Welling, M. 2017. Visualizing deep neural network decisions: Prediction difference analysis. *CoRR* abs/1702.04595.

6 Supplemental material

Other comparisons

We report comparison of our approach (*without preprocessing*) with different uncertainty estimation metrics including entropy (Shannon 1948), margin(Scheffer, Decomain, and Wrobel 2001), MC-dropout (Gal and Ghahramani 2016), and mutual information between predictions and model posterior (Gal, Islam, and Ghahramani 2017). To calculate uncertainty metrics based on MC-dropout we run 100 times the model for each input while all dropout layers are enabled. The results are shown in Tables 3 - 10.

We also compare our results with ODIN and MD (Mahalanobis Distance) *without preprocessing*. Their results are also listed in tables 11-14. Our approach still outperforms these approaches. As expected, the performance of these approaches is reduced because they do not have access to OOD samples anymore.

Hyper parameters

We used three different models for the experiments. Their hyper parameters are listed in table 15. The best perturbation magnitude for each dataset was obtained from the following vector [0.0, 0.0005, 0.001, 0.0015, 0.002, 0.0025, 0.005, 0.01, 0.05, 0.1, 0.15, 0.2]. This vector is fixed for all approaches including ours.

Table 3: Comparison between our approach (without preprocessing) and entropy for the MNIST datasets. Values are in percentages and ↓ indicates that lower values are better, while ↑ indicates that higher values are better. **Bold** indicates the best score.

ID model	OOD	FPR at 95% TPR ↓	Detection error ↓	AUROC ↑	AUPR Out ↑	AUPR In ↑
entropy / ours without preprocessing						
MNIST CUSTOM-CNN	F-MNIST	6.66/ 0.42	5.73/ 2.71	98.29/ 99.04	98.57/ 98.78	98.02/ 99.27
	Omniglot	6.67/ 0.0	5.69/ 0.0	97.94/ 100.0	98.45/ 100.0	97.25/ 100.0
	Gaussian	0.84/ 0.0	0.99/ 0.0	99.72/ 100.0	99.81/ 100.0	99.54/ 100.0
	Uniform	3.38/ 0.0	2.72/ 0.0	98.69/ 100.0	99.11/ 100.0	97.94/ 100.0

Table 4: Comparison between our approach (without preprocessing) and margin for the MNIST datasets. Values are in percentages and ↓ indicates that lower values are better, while ↑ indicates that higher values are better. **Bold** indicates the best score.

ID model	OOD	FPR at 95% TPR ↓	Detection error ↓	AUROC ↑	AUPR Out ↑	AUPR In ↑
margin / ours without preprocessing						
MNIST CUSTOM-CNN	F-MNIST	7.06/ 0.42	5.92/ 2.71	97.88/ 99.04	98.35/ 98.78	96.94/ 99.27
	Omniglot	7.33/ 0.0	6.02/ 0.0	97.43/ 100.0	98.18/ 100.0	95.8/ 100.0
	Gaussian	2.11/ 0.0	1.52/ 0.0	98.48/ 100.0	99.17/ 100.0	94.64/ 100.0
	Uniform	4.2/ 0.0	3.04/ 0.0	97.78/ 100.0	98.67/ 100.0	93.98/ 100.0

Table 5: Comparison between our approach (without preprocessing) and MC-dropout for the MNIST datasets. Values are in percentages and ↓ indicates that lower values are better, while ↑ indicates that higher values are better. **Bold** indicates the best score.

ID model	OOD	FPR at 95% TPR ↓	Detection error ↓	AUROC ↑	AUPR Out ↑	AUPR In ↑
MC-dropout / ours without preprocessing						
MNIST CUSTOM-CNN	F-MNIST	6.59/ 0.42	5.53/ 2.71	98.19/ 99.04	98.58/ 98.78	97.71/ 99.27
	Omniglot	6.54/ 0.0	5.3/ 0.0	97.78/ 100.0	98.41/ 100.0	96.72/ 100.0
	Gaussian	2.3/ 0.0	1.66/ 0.0	98.71/ 100.0	99.24/ 100.0	96.44/ 100.0
	Uniform	5.03/ 0.0	3.71/ 0.0	97.45/ 100.0	98.43/ 100.0	94.43/ 100.0

Table 6: Comparison between our approach (without preprocessing) and mutual information for the MNIST datasets. Values are in percentages and ↓ indicates that lower values are better, while ↑ indicates that higher values are better. **Bold** indicates the best score.

ID model	OOD	FPR at 95% TPR ↓	Detection error ↓	AUROC ↑	AUPR Out ↑	AUPR In ↑
mutual information / ours without preprocessing						
MNIST CUSTOM-CNN	F-MNIST	7.27/ 0.42	5.71/ 2.71	97.53/ 99.04	98.21/ 98.78	96.24/ 99.27
	Omniglot	6.6/ 0.0	5.42/ 0.0	97.6/ 100.0	98.33/ 100.0	95.75/ 100.0
	Gaussian	4.29/ 0.0	2.71/ 0.0	96.72/ 100.0	98.22/ 100.0	89.69/ 100.0
	Uniform	7.1/ 0.0	4.79/ 0.0	95.19/ 100.0	97.3/ 100.0	86.65/ 100.0

Table 7: Comparison between our approach (without preprocessing) and entropy for the CIFAR-10 and CIFAR-100 datasets. ↓ indicates that lower values are better, whereas ↑ indicates that higher values are better. **Bold** indicates the best score.

ID model	OOD	FPR at 95% TPR ↓	Detection error ↓	AUROC ↑	AUPR Out ↑	AUPR In ↑
entropy / ours without preprocessing						
CIFAR10 VGG16	TinyImagenet	35.44/ 25.6	20.22/ 15.29	88.11/ 95.91	90.11/ 95.21	85.16/ 96.52
	LSUN	26.87/ 8.75	15.92/ 6.87	90.8/ 98.21	92.92/ 98.1	88.2/ 98.36
	iSUN	28.19/ 10.66	16.56/ 7.82	90.18/ 98.04	92.49/ 97.93	87.33/ 98.21
	SVHN	27.86/ 8.46	16.42/ 6.7	89.42/ 97.32	92.22/ 97.95	84.6/ 95.64
	Gaussian	20.43/ 0.0	12.24/ 0.0	86.96/ 100.0	92.19/ 100.0	74.06/ 100.0
	Uniform	25.15/ 0.0	14.74/ 0.0	82.25/ 100.0	89.56/ 100.0	67.39/ 100.0
CIFAR10 ResNet-V1-44	TinyImagenet	36.14/ 7.93	20.55/ 6.45	88.28/ 98.32	90.41/ 97.83	85.55/ 98.65
	LSUN	27.58/ 2.49	16.28/ 3.73	91.09/ 99.14	92.97/ 98.99	88.88/ 99.29
	iSUN	29.62/ 7.38	17.31/ 6.19	90.4/ 98.57	92.37/ 98.31	88.18/ 98.81
	SVHN	21.0/ 4.44	13.0/ 4.66	93.77/ 98.68	94.93/ 98.78	92.34/ 98.12
	Gaussian	38.29/ 0.0	21.64/ 0.0	83.47/ 100.0	88.02/ 100.0	72.79/ 100.0
	Uniform	18.39/ 0.0	11.63/ 0.0	93.15/ 100.0	95.05/ 100.0	89.01/ 100.0
CIFAR100 VGG16	TinyImagenet	63.07/ 51.87	34.03/ 28.43	75.36/ 91.14	78.01/ 89.0	71.35/ 92.73
	LSUN	61.31/ 28.14	33.11/ 16.55	74.84/ 95.36	78.63/ 94.62	69.7/ 96.04
	iSUN	63.93/ 31.27	34.46/ 18.13	73.62/ 94.66	77.33/ 93.98	68.48/ 95.42
	SVHN	64.82/ 22.23	34.9/ 13.6	73.32/ 92.5	77.1/ 94.23	69.83/ 87.44
	Gaussian	94.82/ 0.0	48.5/ 0.0	12.55/ 100.0	36.75/ 100.0	32.93/ 100.0
	Uniform	96.49/ 0.0	48.94/ 0.0	8.23/ 100.0	34.52/ 100.0	32.11/ 100.0
CIFAR100 ResNet-V1-44	TinyImagenet	62.59/ 20.45	33.76/ 12.7	77.38/ 96.77	79.72/ 96.32	74.01/ 97.24
	LSUN	61.12/ 19.68	33.05/ 12.33	78.0/ 96.81	80.45/ 96.55	74.76/ 97.21
	iSUN	61.12/ 21.92	33.03/ 13.46	76.95/ 96.37	79.82/ 96.14	73.29/ 96.81
	SVHN	57.17/ 12.15	31.06/ 8.55	81.87/ 96.95	83.74/ 97.4	79.93/ 95.2
	Gaussian	68.39/ 0.0	36.58/ 0.0	50.01/ 100.0	65.5/ 100.0	44.76/ 100.0
	Uniform	47.66/ 0.0	26.32/ 0.0	71.05/ 100.0	80.73/ 100.0	57.96/ 100.0

Table 8: Comparison between our approach (without preprocessing) and margin for the CIFAR-10 and CIFAR-100 datasets. ↓ indicates that lower values are better, whereas ↑ indicates that higher values are better. **Bold** indicates the best score.

ID model	OOD	FPR at 95% TPR ↓	Detection error ↓	AUROC ↑	AUPR Out ↑	AUPR In ↑
margin / ours without preprocessing						
CIFAR10 VGG16	TinyImagenet	36.17/ 25.6	20.54/ 15.29	87.26/ 95.91	89.63/ 95.21	82.51/ 96.52
	LSUN	28.73/ 8.75	16.8/ 6.87	89.71/ 98.21	92.3/ 98.1	85.02/ 98.36
	iSUN	29.94/ 10.66	17.46/ 7.82	89.12/ 98.04	91.88/ 97.93	84.1/ 98.21
	SVHN	29.33/ 8.46	17.14/ 6.7	88.75/ 97.32	91.77/ 97.95	82.97/ 95.64
	Gaussian	21.03/ 0.0	12.54/ 0.0	86.49/ 100.0	91.91/ 100.0	73.3/ 100.0
	Uniform	25.4/ 0.0	14.79/ 0.0	82.08/ 100.0	89.47/ 100.0	67.19/ 100.0
CIFAR10 ResNet-V1-44	TinyImagenet	36.37/ 7.93	20.67/ 6.45	87.3/ 98.32	89.94/ 97.83	82.56/ 98.65
	LSUN	28.37/ 2.49	16.68/ 3.73	89.83/ 99.14	92.35/ 98.99	85.42/ 99.29
	iSUN	30.39/ 7.38	17.69/ 6.19	89.13/ 98.57	91.75/ 98.31	84.59/ 98.81
	SVHN	21.27/ 4.44	13.12/ 4.66	92.51/ 98.68	94.31/ 98.78	89.09/ 98.12
	Gaussian	37.84/ 0.0	21.42/ 0.0	84.6/ 100.0	88.5/ 100.0	77.67/ 100.0
	Uniform	19.1/ 0.0	12.03/ 0.0	92.34/ 100.0	94.62/ 100.0	88.01/ 100.0
CIFAR100 VGG16	TinyImagenet	64.32/ 51.87	34.64/ 28.43	73.21/ 91.14	76.74/ 89.0	67.55/ 92.73
	LSUN	62.63/ 28.14	33.8/ 16.55	73.08/ 95.36	77.53/ 94.62	66.99/ 96.04
	iSUN	65.05/ 31.27	35.03/ 18.13	71.78/ 94.66	76.17/ 93.98	65.8/ 95.42
	SVHN	66.64/ 22.23	35.81/ 13.6	70.95/ 92.5	75.6/ 94.23	65.51/ 87.44
	Gaussian	95.49/ 0.0	48.69/ 0.0	12.09/ 100.0	36.14/ 100.0	32.81/ 100.0
	Uniform	96.79/ 0.0	49.05/ 0.0	7.56/ 100.0	34.12/ 100.0	31.98/ 100.0
CIFAR100 ResNet-V1-44	TinyImagenet	63.65/ 20.45	34.32/ 12.7	73.48/ 96.77	77.47/ 96.32	67.86/ 97.24
	LSUN	62.37/ 19.68	33.67/ 12.33	73.59/ 96.81	77.92/ 96.55	67.52/ 97.21
	iSUN	62.31/ 21.92	33.65/ 13.46	73.08/ 96.37	77.62/ 96.14	67.07/ 96.81
	SVHN	58.79/ 12.15	31.88/ 8.55	76.25/ 96.95	80.44/ 97.4	70.57/ 95.2
	Gaussian	68.77/ 0.0	36.76/ 0.0	51.88/ 100.0	65.89/ 100.0	46.11/ 100.0
	Uniform	48.27/ 0.0	26.54/ 0.0	73.29/ 100.0	81.27/ 100.0	63.15/ 100.0

Table 9: Comparison between our approach (without preprocessing) and MC-dropout for the CIFAR-10 and CIFAR-100 datasets. ↓ indicates that lower values are better, whereas ↑ indicates that higher values are better. **Bold** indicates the best score. We do not report results for ResNet because it does not use dropout layers.

ID model	OOD	FPR at 95% TPR ↓	Detection error ↓	AUROC ↑	AUPR Out ↑	AUPR In ↑
margin / ours without preprocessing						
CIFAR10 VGG16	TinyImagenet	42.38/ 25.6	23.62/ 15.29	87.17/ 95.91	83.77/ 95.21	85.54/ 96.52
	LSUN	26.77/ 8.75	15.87/ 6.87	91.9/ 98.21	92.24/ 98.1	90.24/ 98.36
	iSUN	29.64/ 10.66	17.31/ 7.82	90.69/ 98.04	90.8/ 97.93	88.98/ 98.21
	SVHN	26.39/ 8.46	15.69/ 6.7	91.76/ 97.32	92.65/ 97.95	89.96/ 95.64
	Gaussian	6.49/ 0.0	5.56/ 0.0	98.75/ 100.0	98.94/ 100.0	98.59/ 100.0
	Uniform	6.36/ 0.0	5.43/ 0.0	98.74/ 100.0	98.94/ 100.0	98.58/ 100.0
CIFAR100 VGG16	TinyImagenet	61.02/ 51.87	33.0/ 28.43	78.28/ 91.14	79.9/ 89.0	74.57/ 92.73
	LSUN	45.08/ 28.14	25.04/ 16.55	86.26/ 95.36	87.91/ 94.62	84.19/ 96.04
	iSUN	55.92/ 31.27	30.44/ 18.13	82.63/ 94.66	83.93/ 93.98	80.58/ 95.42
	SVHN	41.88/ 22.23	23.43/ 13.6	87.47/ 92.5	89.25/ 94.23	85.27/ 87.44
	Gaussian	6.37/ 0.0	5.67/ 0.0	98.78/ 100.0	98.89/ 100.0	98.73/ 100.0
	Uniform	5.91/ 0.0	5.43/ 0.0	98.86/ 100.0	98.97/ 100.0	98.81/ 100.0

Table 10: Comparison between our approach (without preprocessing) and mutual information for the CIFAR-10 and CIFAR-100 datasets. \downarrow indicates that lower values are better, whereas \uparrow indicates that higher values are better. **Bold** indicates the best score. We do not report results for ResNet because it does not use dropout layers.

ID model	OOD	FPR at 95% TPR \downarrow	Detection error \downarrow	AUROC \uparrow	AUPR Out \uparrow	AUPR In \uparrow
mutual information / ours without preprocessing						
CIFAR10	TinyImagenet	39.97/ 25.6	22.48/ 15.29	90.13/ 95.91	86.51/ 95.21	90.67/ 96.52
	LSUN	23.29/ 8.75	14.14/ 6.87	95.41/ 98.21	94.93/ 98.1	95.64/ 98.36
	VGG16	27.7/ 10.66	16.34/ 7.82	93.86/ 98.04	93.31/ 97.93	94.04/ 98.21
	SVHN	23.59/ 8.46	14.28/ 6.7	95.1/ 97.32	95.36/ 97.95	95.12/ 95.64
	Gaussian	0.04/ 0.0	1.05/ 0.0	99.94/ 100.0	99.94/ 100.0	99.95/ 100.0
	Uniform	0.03/ 0.0	0.96/ 0.0	99.95/ 100.0	99.95/ 100.0	99.95/ 100.0
CIFAR100	TinyImagenet	58.42/ 51.87	31.68/ 28.43	83.61/ 91.14	83.52/ 89.0	83.38/ 92.73
	LSUN	34.39/ 28.14	19.69/ 16.55	93.19/ 95.36	93.0/ 94.62	93.7/ 96.04
	VGG16	50.75/ 31.27	27.87/ 18.13	89.0/ 94.66	88.41/ 93.98	89.79/ 95.42
	SVHN	27.96/ 22.23	16.47/ 13.6	94.7/ 92.5	94.64/ 94.23	95.05/ 87.44
	Gaussian	0.0/ 0.0	0.29/ 0.0	99.99/ 100.0	99.99/ 100.0	99.99/ 100.0
	Uniform	0.0/ 0.0	0.26/ 0.0	99.99/ 100.0	99.99/ 100.0	99.99/ 100.0

Table 11: Comparison between our approach (without preprocessing) and ODIN (without preprocessing) information for the MNIST datasets. Values are in percentages and \downarrow indicates that lower values are better, while \uparrow indicates that higher values are better. **Bold** indicates the best score.

ID model	OOD	FPR at 95% TPR \downarrow	Detection error \downarrow	AUROC \uparrow	AUPR Out \uparrow	AUPR In \uparrow
ODIN / ours (both without preprocessing)						
MNIST CUSTOM-CNN	F-MNIST	5.75/ 0.42	5.32/ 2.71	98.77/ 99.04	98.85/ 98.78	98.71/ 99.27
	Omniglot	5.33/ 0.0	4.84/ 0.0	98.63/ 100.0	98.92/ 100.0	98.23/ 100.0
	Gaussian	0.13/ 0.0	0.15/ 0.0	99.96/ 100.0	99.97/ 100.0	99.91/ 100.0
	Uniform	0.8/ 0.0	1.04/ 0.0	99.74/ 100.0	99.81/ 100.0	99.58/ 100.0

Table 12: Comparison between our approach (without preprocessing) and MD (without preprocessing) information for the MNIST datasets. Values are in percentages and \downarrow indicates that lower values are better, while \uparrow indicates that higher values are better. **Bold** indicates the best score.

ID model	OOD	FPR at 95% TPR \downarrow	Detection error \downarrow	AUROC \uparrow	AUPR Out \uparrow	AUPR In \uparrow
MD / ours (both without preprocessing)						
MNIST CUSTOM-CNN	F-MNIST	46.43/ 0.42	25.69/ 2.71	82.71/ 99.04	86.09/ 98.78	78.39/ 99.27
	Omniglot	27.59/ 0.0	16.27/ 0.0	91.7/ 100.0	93.42/ 100.0	89.9/ 100.0
	Gaussian	41.57/ 0.0	22.78/ 0.0	66.29/ 100.0	79.97/ 100.0	53.53/ 100.0
	Uniform	42.2/ 0.0	23.56/ 0.0	72.13/ 100.0	82.32/ 100.0	58.33/ 100.0

Table 13: Comparison between our approach (without preprocessing) and ODIN (without preprocessing) for the CIFAR-10 and CIFAR-100 datasets. \downarrow indicates that lower values are better, whereas \uparrow indicates that higher values are better. **Bold** indicates the best score.

ID model	OOD	FPR at 95% TPR \downarrow	Detection error \downarrow	AUROC \uparrow	AUPR Out \uparrow	AUPR In \uparrow
ODIN / ours (both without preprocessing)						
CIFAR10 VGG16	TinyImagenet	33.18/ 25.6	19.08/ 15.29	90.91/ 95.91	91.58/ 95.21	89.6/ 96.52
	LSUN	20.75/ 8.75	12.86/ 6.87	94.42/ 98.21	95.38/ 98.1	93.17/ 98.36
	iSUN	22.5/ 10.66	13.74/ 7.82	93.77/ 98.04	94.82/ 97.93	92.45/ 98.21
	SVHN	23.95/ 8.46	14.46/ 6.7	92.07/ 97.32	93.82/ 97.95	88.7/ 95.64
	Gaussian	26.95/ 0.0	15.9/ 0.0	85.36/ 100.0	90.61/ 100.0	73.12/ 100.0
	Uniform	50.5/ 0.0	27.74/ 0.0	70.1/ 100.0	79.7/ 100.0	57.25/ 100.0
CIFAR10 ResNet-V1-44	TinyImagenet	29.75/ 7.93	17.37/ 6.45	91.67/ 98.32	92.89/ 97.83	89.61/ 98.65
	LSUN	17.43/ 2.49	11.2/ 3.73	95.19/ 99.14	96.14/ 98.99	93.6/ 99.29
	iSUN	19.0/ 7.38	12.0/ 6.19	94.67/ 98.57	95.7/ 98.31	93.08/ 98.81
	SVHN	15.12/ 4.44	10.06/ 4.66	96.25/ 98.68	96.8/ 98.78	95.43/ 98.12
	Gaussian	85.7/ 0.0	45.32/ 0.0	34.42/ 100.0	50.86/ 100.0	38.68/ 100.0
	Uniform	38.61/ 0.0	21.78/ 0.0	77.72/ 100.0	85.46/ 100.0	63.84/ 100.0
CIFAR100 VGG16	TinyImagenet	53.49/ 51.87	29.24/ 28.43	81.29/ 91.14	83.73/ 89.0	76.85/ 92.73
	LSUN	49.32/ 28.14	27.16/ 16.55	81.47/ 95.36	84.74/ 94.62	75.92/ 96.04
	iSUN	53.33/ 31.27	29.17/ 18.13	80.55/ 94.66	83.75/ 93.98	74.85/ 95.42
	SVHN	49.3/ 22.23	27.13/ 13.6	82.31/ 92.5	85.3/ 94.23	77.81/ 87.44
	Gaussian	81.97/ 0.0	42.85/ 0.0	27.22/ 100.0	50.36/ 100.0	36.68/ 100.0
	Uniform	82.37/ 0.0	43.09/ 0.0	26.78/ 100.0	50.06/ 100.0	36.57/ 100.0
CIFAR100 ResNet-V1-44	TinyImagenet	52.91/ 20.45	28.94/ 12.7	82.26/ 96.77	84.43/ 96.32	78.31/ 97.24
	LSUN	42.07/ 19.68	23.53/ 12.33	86.44/ 96.81	88.6/ 96.55	82.7/ 97.21
	iSUN	46.58/ 21.92	25.78/ 13.46	84.32/ 96.37	86.87/ 96.14	80.07/ 96.81
	SVHN	34.33/ 12.15	19.64/ 8.55	88.98/ 96.95	91.14/ 97.4	84.98/ 95.2
	Gaussian	98.71/ 0.0	49.69/ 0.0	2.83/ 100.0	31.77/ 100.0	31.16/ 100.0
	Uniform	88.34/ 0.0	45.86/ 0.0	19.23/ 100.0	43.7/ 100.0	34.6/ 100.0

Table 14: Comparison between our approach (without preprocessing) and MD (without preprocessing) for the CIFAR-10 and CIFAR-100 datasets. \downarrow indicates that lower values are better, whereas \uparrow indicates that higher values are better. **Bold** indicates the best score.

ID model	OOD	FPR at 95% TPR \downarrow	Detection error \downarrow	AUROC \uparrow	AUPR Out \uparrow	AUPR In \uparrow
MD / ours (both without preprocessing)						
CIFAR10 VGG16	TinyImagenet	31.94/ 25.6	18.47/ 15.29	90.46/ 95.91	91.42/ 95.21	88.68/ 96.52
	LSUN	21.46/ 8.75	13.21/ 6.87	93.66/ 98.21	94.77/ 98.1	92.27/ 98.36
	iSUN	23.71/ 10.66	14.34/ 7.82	92.87/ 98.04	94.17/ 97.93	91.32/ 98.21
	SVHN	33.48/ 8.46	19.22/ 6.7	87.57/ 97.32	90.4/ 97.95	79.89/ 95.64
	Gaussian	13.65/ 0.0	8.89/ 0.0	93.49/ 100.0	95.75/ 100.0	89.29/ 100.0
	Uniform	19.05/ 0.0	11.51/ 0.0	87.98/ 100.0	92.79/ 100.0	75.67/ 100.0
CIFAR10 ResNet-V1-44	TinyImagenet	81.25/ 7.93	43.08/ 6.45	61.35/ 98.32	64.74/ 97.83	56.26/ 98.65
	LSUN	80.59/ 2.49	42.78/ 3.73	56.82/ 99.14	62.99/ 98.99	50.68/ 99.29
	iSUN	82.21/ 7.38	43.6/ 6.19	56.69/ 98.57	62.04/ 98.31	51.19/ 98.81
	SVHN	65.94/ 4.44	35.47/ 4.66	76.25/ 98.68	78.12/ 98.78	71.09/ 98.12
	Gaussian	0.0/ 0.0	0.0/ 0.0	100.0/ 100.0	100.0/ 100.0	100.0/ 100.0
	Uniform	0.0/ 0.0	0.01/ 0.0	100.0/ 100.0	100.0/ 100.0	100.0/ 100.0
CIFAR100 VGG16	TinyImagenet	58.21/ 51.87	31.6/ 28.43	79.09/ 91.14	81.31/ 89.0	74.68/ 92.73
	LSUN	55.22/ 28.14	30.06/ 16.55	80.08/ 95.36	82.59/ 94.62	75.71/ 96.04
	iSUN	58.72/ 31.27	31.86/ 18.13	78.86/ 94.66	81.16/ 93.98	74.17/ 95.42
	SVHN	53.26/ 22.23	29.11/ 13.6	77.03/ 92.5	81.99/ 94.23	69.74/ 87.44
	Gaussian	51.73/ 0.0	27.84/ 0.0	58.68/ 100.0	74.3/ 100.0	48.91/ 100.0
	Uniform	67.74/ 0.0	36.14/ 0.0	42.55/ 100.0	62.8/ 100.0	41.71/ 100.0
CIFAR100 ResNet-V1-44	TinyImagenet	86.99/ 20.45	45.99/ 12.7	56.66/ 96.77	59.31/ 96.32	52.74/ 97.24
	LSUN	87.16/ 19.68	46.06/ 12.33	54.02/ 96.81	58.24/ 96.55	50.01/ 97.21
	iSUN	90.77/ 21.92	47.88/ 13.46	50.76/ 96.37	53.47/ 96.14	48.17/ 96.81
	SVHN	81.16/ 12.15	43.07/ 8.55	62.49/ 96.95	65.85/ 97.4	57.26/ 95.2
	Gaussian	0.0/ 0.0	0.0/ 0.0	100.0/ 100.0	100.0/ 100.0	100.0/ 100.0
	Uniform	0.04/ 0.0	0.14/ 0.0	99.99/ 100.0	99.99/ 100.0	99.97/ 100.0

Table 15: Hyper parameters used for different models in the experiments

Model	Optimizer	Epochs	Batch size	learning rate	Dropout layer	Data augmentation
ResNet	Adam	200	32	0.001	No	Yes
VGG-16	SGD	250	128	0.1	Yes	Yes
CNN (MNIST)	Adam	100	5000	0.001	Yes	Yes