

Understanding the Feature Norm for Out-of-Distribution Detection

Jaewoo Park^{1,2} Jacky Chen Long Chai¹ Jaeho Yoon¹ Andrew Beng Jin Teoh^{1†}
¹Yonsei University ²AiV Co.

Abstract

A neural network trained on a classification dataset often exhibits a higher vector norm of hidden layer features for in-distribution (ID) samples, while producing relatively lower norm values on unseen instances from out-of-distribution (OOD). Despite this intriguing phenomenon being utilized in many applications, the underlying cause has not been thoroughly investigated. In this study, we demystify this very phenomenon by scrutinizing the discriminative structures concealed in the intermediate layers of a neural network. Our analysis leads to the following discoveries: (1) The feature norm is a confidence value of a classifier hidden in the network layer, specifically its maximum logit. Hence, the feature norm distinguishes OOD from ID in the same manner that a classifier confidence does. (2) The feature norm is class-agnostic, thus it can detect OOD samples across diverse discriminative models. (3) The conventional feature norm fails to capture the deactivation tendency of hidden layer neurons, which may lead to misidentification of ID samples as OOD instances. To resolve this drawback, we propose a novel negative-aware norm (NAN) that can capture both the activation and deactivation tendencies of hidden layer neurons. We conduct extensive experiments on NAN, demonstrating its efficacy and compatibility with existing OOD detectors, as well as its capability in label-free environments.

1. Introduction

Deep learning-based models are increasingly used for safety-critical applications such as autonomous driving and medical diagnosis. Despite the effectiveness of deep models in closed-set environments where all test queries are sampled from the same distribution of train data, the deep models are reported fairly vulnerable [33, 16] to outliers from out-of-distribution [19, 51] and make highly confident but invalid predictions thereon [35]. As it is critical to prevent such malfunction in deploying deep models for open environment applications, the out-of-distribution (OOD) detec-

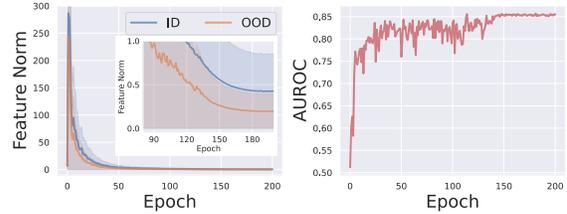


Figure 1: (left) As a discriminative model is trained, its hidden layer features exhibit higher vector norm on in-distribution samples (ID) and relatively lower norm on out-of-distribution (OOD) instances. This phenomenon prevails even when the model reduces the overall feature norm (e.g. by weight decay). (right) As a result, ID samples are separated from OOD instances with respect to the feature norm. To see its underlying cause, we analyze the discriminative structures concealed in the hidden layer.

tion problem has attracted massive attention in recent years [52].

Despite the importance of this field, only a handful of works have been devoted to understanding how the deep network becomes aware of OOD [9, 10, 8, 30, 31]. One particular under-studied signal in OOD detection is *the norm of feature vectors* residing in the hidden layers of neural networks. Its known behavior is that a model trained on the ID data exhibits larger values of feature norm over ID samples than the OOD instances [7, 53, 3, 28]. However, the studies are mainly empirical and provide no underlying principle of the feature norm at a fundamental level.

A preliminary attempt at understanding the feature norm has been given in the appendix of [45]. The authors of [45] argue that minimizing the cross entropy (CE) maximizes the feature norm of ID samples. However, the argument is not general. As we observe in Fig. 1, training the weight-decayed model decreases the overall feature norm, but the separation between ID and OOD remains obvious. Hence, we require a new lens to understand the underlying cause of feature norm separation.

In this work, we study *why* the feature norm separates ID from OOD. To this end, we both theoretically and empirically show that the feature norm is equal to a confidence value of a classifier hidden in the corresponding layer.

[†] Corresponding author: Andrew Beng Jin Teoh

Based on the existing theory on the classifier confidence [10], the equality guarantees the detection capability of feature norm.

Furthermore, our analysis indicates that the feature norm is agnostic to the class label space. This suggests that the feature norm can detect OOD using any general discriminative model, including self-supervised classifiers. We validate this postulation empirically under several aspects: Firstly, by considering inter- and intra-class learning independently, we show that inter-class learning enables the feature norm to separate OOD from the training fold of ID. The intra-class learning, on the other hand, generalizes the detection capability to the test environment, enabling the feature norm to differentiate OOD from the test fold of ID. The finding shows that inter- and intra-class learning corresponds to memorization and generalization, respectively, in the context of OOD detection. Secondly, we show that the detection capability of feature norm is strongly correlated to the entropy of activation (*i.e.* diversity of on/off status of neurons). As activation entropy is a class-agnostic characteristic, the finding reinforces our postulation.

In addition to that, we observe that the conventional vector norm only captures the activation tendency of hidden layer neurons, but misses the deactivation counterpart. Failing to account for the deactivation tendencies results in the loss of important characteristics specific to ID samples, potentially leading to misidentification of such instances as OOD examples. Motivated by this drawback, we derive a novel negative-aware norm that captures both the activation and deactivation tendencies of hidden layer neurons.

We perform a thorough assessment of the NAN and demonstrate its efficacy across OOD benchmarks. Additionally, we confirm that NAN is compatible with several state-of-the-art OOD detectors. Furthermore, NAN is free of hyperparameters, requires no classification layer, and does not necessitate expensive feature extraction from a bank set. Consequently, NAN can be readily deployed in scenarios where class labels are unavailable. We evaluate NAN in unsupervised environments using self-supervised models and assess its performance on one-class classification benchmarks.

The contributions of our work are summarized as follows:

- We demystify the OOD detection capability of the feature norm by showing that the feature norm is a confidence value of a classifier hidden in the corresponding layer (Sec. 3).
- We reveal that the feature norm is class-agnostic, hence able to detect OOD using general discriminative models (Sec. 4). We validate this property under several aspects including inter/intra-class learning and activation entropy.

- We put forward a novel negative-aware norm (NAN), which captures both activation and deactivation tendencies of hidden layer neurons (Sec. 5). NAN is hyperparameter-free, label-free, and bank-set-free. NAN can be easily integrated with state-of-the-art OOD detectors. (Sec. 6)

2. Background

The goal of OOD detection is to devise a score function $S(\mathbf{x})$ that determines an input sample \mathbf{x} as OOD if $S(\mathbf{x}) < \tau$ for some threshold τ and as ID otherwise. There are several ways to derive such a score function from a discriminative model $p_\theta(y|\mathbf{x})$. A standard detection score is the maximum softmax probability (MSP) score [16], which is defined as $S(\mathbf{x}) = \max_y p_\theta(y|\mathbf{x})$ with p_θ modeled by the softmax function.

Other OOD detection scores include the energy score [27] that extracts the energy function [13] from the classification layer. [15] proposes the KL divergence to the uniform prediction, while [45] applies only the maximum value of logit.

Other works propose the utilization of distance metrics for OOD detection. [25] applied the Mahalanobis distance as an OOD detector based on a strong parametric assumption that each ID class follows a Gaussian distribution with a shared covariance. A unified approach SSD [40] generalizes the principle of [25], exploiting class clusters attained by unsupervised K -means. As SSD requires no class labels, its usage is general and applicable to both supervised and unsupervised models. ViM [46] adopts SSD but uses the orthogonal distance from principal components instead, and combines it with the energy score with manual calibration. CSI [43], on the other hand, defines the detection score by combining a rotation classifier with the k -nearest neighbor distance. The effectiveness of CSI, however, comes from a deliberate design of image-specific data augmentations. As a simpler and model-agnostic approach, [42] proposed the k -nearest neighbor (KNN) distance for OOD detection. Despite its broad applicability [29], KNN requires a careful hyperparameter search on the sampling ratio of the ID bank set and the number of neighbors.

Apart from the distance-based OOD detectors, an alternative approach to detecting OOD is by perturbing the signal of the network. [26, 20] observed a particular input perturbation perturbs OOD samples severely but makes ID samples remain mostly invariant. [41] proposed a rectification layer that clips out all values greater than a given threshold. Despite their effectiveness, the perturbation methods rely on specific assumptions of network signal distributions and are sensitive to hyperparameters.

On feature norm. The first application of feature norm for OOD detection was reported by [7], whose authors observed that the magnitude (l_2 -norm) of embedding vector tends to be larger for ID than OOD. The same trend was ob-

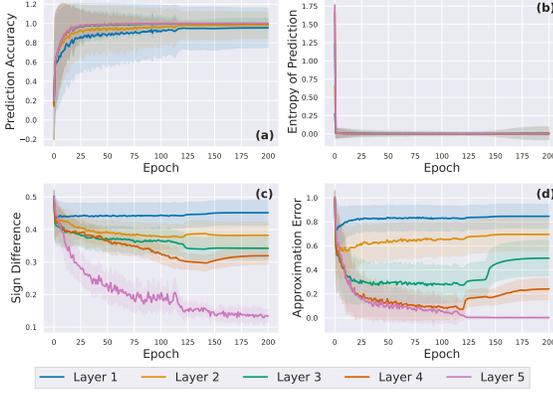


Figure 2: **The results on hidden classifiers** of MLP-5 trained on CIFAR-10 (ID). (a) The prediction accuracy of the hidden classifier increases through learning. (b) Accordingly, the prediction becomes more deterministic (*i.e.*, confident). (c,d) As the sign difference between the feature vector and class weight $\mathbf{c}_y^{(l)}$ is reduced, the approximation error between the feature norm and the maximum value of the hidden classifier is reduced in a similar trend, *verifying* our Thm. 3. Results with other activation functions are in Sec. A.3.1.

served in the appendixes of [43, 45, 21] for generic images. In biometrics, [53] observed the same phenomenon for face images, thereby devising a score that can more effectively reject unseen identities based on the feature norm. [28] extended the application of feature norm, showing that it can measure the quality score of the face image. On the other hand, [3, 4] observed that the norm of feature embedding effectively differentiates a person from his/her surrounding background, and thus can be used to improve the performance and efficiency of person search. Besides OOD detection, [55] observed that the embedding vectors of highly discriminative samples lie in the area of the large norm. [50] extended this observation, demonstrating the samples with large feature norms are not only more discriminative but also more transferable for domain adaptation.

Although numerous works report empirical observations of the phenomenon, to our best knowledge, no work provides a systematic theoretical explanation of the underlying mechanism of feature norm.

3. Understanding Feature Norm as a Confidence of Hidden Classifier

In this section, we show that the feature norm is a confidence value of a discriminative classifier covertly concealed in the corresponding layer. Specifically, under a regularity condition, the l_1 -norm of the feature vector is equal to the maximum logit of a hidden classifier attained by binarizing the network weights. Hence, based on the theory from [9],

the feature norm is guaranteed its detection proficiency.

3.1. Theoretical analysis

Notation and setup. Let $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ be the train ID dataset where $y_i \in \mathcal{Y} = \{1, \dots, K\}$ are labels from K classes. Suppose our model is a multi-layer perceptron (MLP) whose l -th hidden layer consists of the d_l -dimensional feature vector $\mathbf{a}^{(l)}$ computed by $\mathbf{a}^{(l)} = \sigma(\mathbf{W}^{(l)T} \mathbf{a}^{(l-1)})$ consecutively from the initial layer $l=0$ to the *last hidden layer* $l=L$, where $\mathbf{a}^{(0)} = \mathbf{x}$. The vector of pre-activated units is denoted by $\mathbf{z}^{(l)}$, which satisfies $\mathbf{a}^{(l)} = \sigma(\mathbf{z}^{(l)})$. The activation function σ is assumed to be a unit-wise rectifier such as ReLU [32, 11], GeLU [17], and Leaky ReLU [48]. Each weight matrix $\mathbf{W}^{(l)} \in \mathbb{R}^{d_{l-1} \times d_l}$ constitutes trainable parameters θ . The classifier logit $\psi(\mathbf{x}) \in \mathbb{R}^K$ is computed by $\psi(\mathbf{x}) = \mathbf{W}^{(L+1)T} \mathbf{a}^{(L)}$.

Assumption. We assume arbitrary class type for the label space \mathcal{Y} ; classes can be supervised labels, instance classes, or even noisy labels.

To extract a hidden classifier from each hidden layer of the model, we first access the hidden layer through matrix multiplication.

Proposition 1. *The final logit is represented by*

$$\psi(\mathbf{x}) = \mathbf{C}^{(l)} \mathbf{a}^{(l)} \quad (1)$$

for each hidden layer l , where

$$\mathbf{C}^{(l)} = \left(\prod_{k=0}^{L-l-1} \mathbf{W}^{(L+1-k)T} \mathbf{D}^{(L-k)} \right) \mathbf{W}^{(l+1)T} \quad (2)$$

with $\mathbf{D}^{(l)} = \text{diag}(\frac{\sigma(z_1)}{z_1}, \dots, \frac{\sigma(z_{d_l})}{z_{d_l}})$ and the convention $\frac{\cdot}{0} = 0$. The matrix $\mathbf{C}^{(l)} = \mathbf{C}^{(l)}(\mathbf{x}) \in \mathbb{R}^{K \times d_l}$ depends on \mathbf{x} .

Proof. All proofs are given in Sec. A. \square

The multiplication by the coefficient matrix $\mathbf{C}^{(l)} = [\mathbf{c}_1^{(l)}, \dots, \mathbf{c}_K^{(l)}]^T$ resembles a classification layer with the column weight $\mathbf{c}_k^{(l)} = \mathbf{c}_k^{(l)}(\mathbf{x})$ as the k -th class proxy.

We note that ψ is called a *discriminative classifier* since the target class unit of logit is maximum $\psi_y > \psi_k$ for all $k \neq y$. If the output classifier ψ is sufficiently discriminative, then binarizing the coefficient matrix $\mathbf{C}^{(l)}$ does not alter the prediction of the classifier. This leads us to a *hidden classifier* $\bar{\psi}^{(l)} \in \mathbb{R}^K$ defined by binarizing the network weights:

$$\bar{\psi}^{(l)}(\mathbf{x}) := \mathbf{B}^{(l)} \mathbf{a}^{(l)} := \text{sign}(\mathbf{C}^{(l)}) \mathbf{a}^{(l)} \quad (3)$$

where $\text{sign}(x) = 1$ if $x > 0$ and -1 otherwise.

Proposition 2. *For all labeled sample (\mathbf{x}, y) , suppose the discriminative learning of $\psi_k(\mathbf{x}) = \mathbf{c}_k^{(l)} \cdot \mathbf{a}^{(l)}$ increases and*

decreases the cosine similarities between $\mathbf{c}_k^{(l)}$ and $\mathbf{a}^{(l)}$ sufficiently for $k=y$ and $k \neq y$, respectively. Then $\bar{\psi}^{(l)}$ is a discriminative classifier with $\bar{\psi}_y^{(l)} > \bar{\psi}_k^{(l)}$ for all $k \neq y$.

In the sufficient condition of Prop. 2, the network aligns the activation pattern $\text{sign}(\mathbf{a}^{(l)})$ [14] with the binary weight $\mathbf{b}_y^{(l)}$ that corresponds to the target class y . Here, $\mathbf{b}_y^{(l)}$ is the y -th row of $\mathbf{B}^{(l)}$. Due to the alignment, the feature norm becomes the prediction confidence $\max_k \bar{\psi}_k^{(l)}(\mathbf{x})$ of the hidden classifier.

Theorem 3. *Given the sufficient condition of Proposition 2, the feature norm*

$$\|\mathbf{a}^{(l)}\|_1 \text{ converges to } \bar{\psi}_y^{(l)}(\mathbf{x}) = \max_k \bar{\psi}_k^{(l)}(\mathbf{x}), \quad (4)$$

in which case $\text{sign}(\mathbf{a}^{(l)}) = \mathbf{b}_y^{(l)}$. In general, for any k

$$0 \leq \|\mathbf{a}^{(l)}\|_1 - \bar{\psi}_k(\mathbf{x}) \leq \|\mathbf{a}^{(l)}\|_\infty \|\text{sign}(\mathbf{a}^{(l)}) - \mathbf{b}_k^{(l)}\|_1 \quad (5)$$

Existing OOD theories on classifiers [9, 10] assure that OOD samples have smaller prediction confidence than ID under regularity conditions. In this case, the feature norm of OOD also has a smaller value due to Thm. 3:

Corollary 4. *If $\max_k \bar{\psi}_k^{(l)}(\mathbf{x}_{ood})$ is sufficiently small, then $\|\mathbf{a}^{(l)}(\mathbf{x}_{ood})\|_1 < \|\mathbf{a}^{(l)}(\mathbf{x}_{ind})\|_1$ for all ID samples \mathbf{x}_{ind} .*

3.2. Empirical verification

We empirically verify the above claims. We train a 5-layer MLP on CIFAR10 (ID) [24]. The full empirical setup is given in Sec. A.3. Fig. 2 shows that the hidden classifiers learn to increase their prediction accuracy while reducing the prediction uncertainty (entropy), verifying Prop. 2. As described in Thm. 3, the discriminative training induces the sign alignment between the hidden layer feature and corresponding class weight $\mathbf{c}_y^{(l)}$, thereby reducing the gap between the feature norm and the maximum confidence of the hidden classifier.

Remark We remark that the trend of approximation error may not be precisely aligned with that of the sign difference (Fig. 2) as the sign difference is the sufficient condition but not a necessary one. Hence, when the sign difference is large, the approximation error can be either large or small; *i.e.* they can be misaligned. However, due to its sufficiency, if the sign difference converges to 0, then the approximation error also decreases to 0.

4. Class Agnosticity of Feature Norm

The theoretical properties of feature norm proven in Sec. 3 hold true with respect to any type of label space,

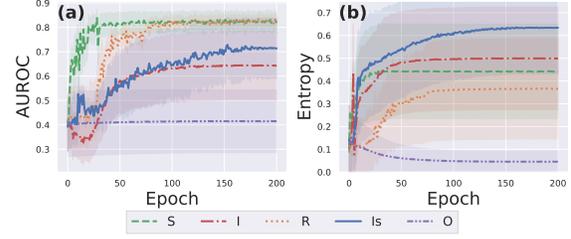


Figure 3: The results on ResNet-18 trained on CIFAR-10 (ID). (a) Training the model increases the OOD detection performance of feature norm if and only if the model is discriminative. (b) Accordingly, training the model increases the entropy of activation if and only if the model is discriminative. Here, the models with S, I, R, and Is labeling schemes are discriminative, while model O is not discriminative.

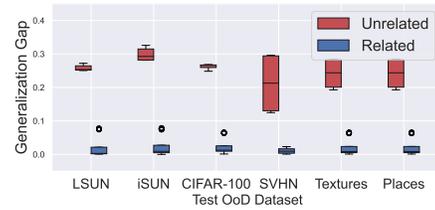


Figure 4: When the intra-class samples are *related semantically* (*i.e.* {S,Is,O}), the OOD detection performance is generalized to test environments (*i.e.* small generalization gap). However, if intra-class samples are randomly related (R), or there is no more than one sample in each class (I), no generalization is observed.

suggesting that the feature norm is class-agnostic and capable of detecting out-of-distribution (OOD) samples with any discriminative model. In this section, we conduct empirical analyses to validate this hypothesis across different aspects. Specifically, we observe that inter/intra-class learning generally enhances the feature norm’s performance. We then demonstrate that the feature norm’s performance is correlated with the entropy of activation, which is another class-agnostic characteristic of neural networks. The feature norm’s dependence on class-agnostic factors provides further evidence supporting our hypothesis.

4.1. Impact of inter/intra-class learning

Setup. We train a ResNet-18 on CIFAR-10, and test against different OODs, *i.e.*, LSUN [54], iSUN [49], CIFAR-100 [24], SVHN [34], Texture [6], and Places [56].

We consider five different training schemes by varying the label space. ‘S’: the supervised learning with generic object categories. ‘I’: the instance-discrimination learning with $y_i=i$. ‘Is’: instance-discrimination with data augmentation (*i.e.* conventional self-supervision). ‘R’: learning with random binary labels. ‘O’: non-discriminative learn-

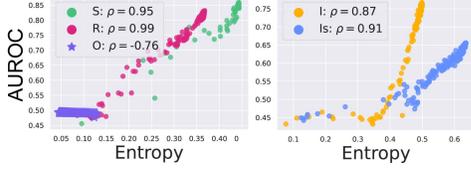


Figure 5: For discriminative models $\{S,R,I,Is\}$, the OOD detection performance of feature norm is positively *correlated* to the averaged entropy of activation (Eq. (7)). However, no consistent correlation is found in the non-discriminative model O.

ing with every ID sample labeled by the same label ‘1’.

The detection score we use is the feature norm $\|\mathbf{a}^{(L)}\|_1$ of the last hidden layer feature $\mathbf{a}^{(L)}$. The performance is measured by the area under receiving operating characteristic curve (AUROC). A more detailed description of the setup and full experimental results are given in Sec. B.

Inter-class learning. To analyze the effect of inter-class learning, we divide the training schemes into two: discriminative learning $\{S,R,I,Is\}$, and non-discriminative learning $\{O\}$. Fig. 3 demonstrates that the feature norm separates OOD from the train fold of ID if and only if the model is trained with inter-class learning. In particular, the feature can detect OOD even if the model is trained with random noisy labels, indicating that its detection capability is independent of the class type of label space.

Intra-class learning. To examine the impact of intra-class learning, we divide the training schemes into two groups $\{S,Is,O\}$ and $\{R,I\}$. In the former group $\{S,Is,O\}$, the intra-class samples are semantically related. On the latter group $\{R,I\}$, there is no semantic relation within the intra-class samples. Fig. 4 indicates the generalization gap between train and test performances for OOD detection. The results support that the detection capability of feature norm is generalized to the test environment if and only if the intra-class samples are semantically related.

Summary on inter/intra-class learning. The detection capability of feature norms does not depend on a particular type of class label. Instead, any type of inter-class learning allows the feature norm to differentiate OOD from the training fold of ID. On the other hand, intra-class learning with any appropriate semantics facilitates the separation of OOD from the test fold of ID. In general, inter-class learning corresponds to memorization, while intra-class is associated with generalization.

4.2. The relation to the entropy of activation

The feature norm’s detection capability depends on the model’s discriminative nature, not the class type. Here, we further show that the capability relies on the entropy of activation, which is another class-agnostic characteristic.

If the model is discriminative, target logits $\bar{\psi}_y^{(L)}(\mathbf{x})$ with

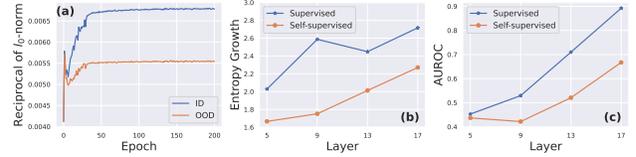


Figure 6: (a) The sparsity of activations, measured by $\|\mathbf{a}^{(L)}\|_0^{-1}$, is maximized and higher on ID samples than on OOD instances. (b) The entropy growth is larger in deeper layers. (c) The OOD detection performance is accordingly better in deeper layers.

different y is maximized for ID samples x . Then, due to

$$\bar{\psi}_y^{(L)}(\mathbf{x}) = \sum_i b_{y,i}^{(L)} a_i^{(L)} \quad (6)$$

with $\mathbf{b}_y^{(L)} = (b_{y,1}^{(L)}, \dots, b_{y,d_L}^{(L)}) \in \mathbb{R}^{d_L}$, the unit $a_i^{(L)}$ is maximized for samples \mathbf{x} in $\{\mathbf{x} : b_{y,i}^{(L)} = 1\}$, and minimized for \mathbf{x} in $\{\mathbf{x} : b_{y,i}^{(L)} = -1\}$. Consequently, the *entropy of activation* is maximized

$$H(a_i^{(L)}) = - \sum_{c=0}^1 \mathbb{P}(1_{a_i^{(L)} > 0} = c) \log \mathbb{P}(1_{a_i^{(L)} > 0} = c) \quad (7)$$

for each neuron $a_i^{(L)}$ of ID samples.

Conversely, if the model is not discriminative, *i.e.*, $\mathcal{Y} = \{1\}$, then all ID samples likely have the same constant binary indicator; $b_{y,i}^{(L)} = c$ for all samples \mathbf{x} where $c \in \{-1, 1\}$. Hence, the activation entropy is minimized in this case.

This trend is empirically validated in Fig. 3b; only discriminative models maximize the activation entropy. Moreover, demonstrated by the strong correlation depicted in Fig. 5, the detection performance of the feature norm depends on the activation entropy, which is a characteristic independent of the class type of the label space.

5. Method: Negative-Aware Norm (NAN)

A missing component in the conventional norm. The network training tends to maximize the confidence of hidden classifier

$$\max \bar{\psi}_y(\mathbf{x}) = \overbrace{\sum_{i: b_{y,i}^{(L)} = 1} a_i^{(L)}}^{\text{maximized}} - \overbrace{\sum_{j: b_{y,j}^{(L)} = -1} a_j^{(L)}}^{\text{minimized}} \quad (8)$$

on ID samples \mathbf{x} under a regularity condition (Prop. 2). This maximization is stronger on ID samples than on OOD instances [9], and hence serves as a key factor that separates OOD from ID (Cor. 4).

The maximization of confidence can be disentangled to maximization of the positive summand

$A := \sum_{i: b_{y,i}^{(L)} = 1} a_i^{(L)}$ and minimization of the negative summand $D := \sum_{j: b_{y,j}^{(L)} = -1} a_j^{(L)}$, which correspond to activation and deactivation of neurons, respectively. The conventional l_1 feature norm $\|\mathbf{a}\|_1$ captures the maximization trend of activation neurons as the summand A converges to $\|\mathbf{a}\|_1$. However the l_1 norm fails to reflect the deactivation responses as the negative summand is diminished with $D \approx 0$ due to the nature of the activation function (e.g. ReLU). Hence, this can lead to potential misidentification of ID samples when the naive l_1 norm is used for OOD detection.

Derivation. To mitigate this drawback, we capture the *deactivation tendency* by the sparsity of activations $\|\mathbf{a}^{(L)}\|_0^{-1}$. The sparsity term reflects the number of deactivated neurons by

$$\|\mathbf{a}^{(L)}\|_0 = d_L - |\{i : a_i^{(L)} \leq 0\}|. \quad (9)$$

Combining the sparsity term with the conventional vector norm, we derive a novel *negative-aware norm (NAN)*

$$\|\mathbf{a}\|_{\text{NAN}} = \|\mathbf{a}^{(L)}\|_1 \cdot \|\mathbf{a}^{(L)}\|_0^{-1}. \quad (10)$$

NAN captures both the activation and deactivation tendencies of ID samples’ neurons. Fig. 6a shows the sparsity term is higher on ID samples than OOD instances, demonstrating that the deactivation tendency is stronger in ID samples’ neurons. Hence, capturing the deactivation tendency likely improves the conventional norm. We conduct extensive experiments on NAN in the next section to validate its effectiveness.

We remark that similar to the l_1 feature norm, the negative-aware norm (NAN) exhibits class-agnostic characteristics, as verified through analyses of inter/intra-class learning and activation entropy in Sec. B.

Additional consideration. We utilize the last hidden layer $\mathbf{a} = \mathbf{a}^{(L)}$ for OOD detection as the last hidden layer exhibits a higher growth in activation entropy and accordingly better performance (Fig. 6bc).

6. Experiments on NAN

The objective of this experiment is to assess the OOD detection capabilities of NAN across diverse configurations using general discriminative models. To achieve this goal, we evaluate NAN’s performance using both supervised and self-supervised models, and assess it in large-scale and small-scale benchmarks, including the one-class classification setting. Additionally, we consider the compatibility of NAN, namely, whether NAN can be combined with other detectors for performance gain. We conclude this section with ablation studies of NAN. A detailed description of the complete experiment setup can be found in Sec. C.

Evaluation metrics The performance is reported by the widely-used metrics: (1) the area under the receiver operating characteristic curve (AUROC), (2) the false positive

rate (FPR95) on the OOD samples when the true positive rate of ID samples is at 95%, (3) closed-set classification accuracy (ACC) of ID.

6.1. Evaluation on large-scale benchmark

Setup. We utilize a ResNet-50 trained on ImageNet-1k. The model is trained either by (1) supervised labels using the contrastive loss [23] or (2) self-supervised instance discrimination loss using momentum embeddings [5]. In the case of the supervised contrastive learning, the classification layer is learned after training and freezing the backbone representation. For fair comparison, all detection scores are applied on the same backbone.

Following the widely-used ImageNet-1k benchmark [22], we test against four test OOD datasets: fine-grained plant images of iNaturalist [44], scene images from SUN [47] and Places [56], and texture images from Texture [6]. All OOD datasets are processed so that no overlapping category is present with ImageNet-1k.

Results. Table 1 shows that NAN is comparable to the state-of-the-art detectors on the ImageNet-1k benchmark. Compared to the OOD detection scores that require a supervised classification layer (i.e. MSP, Energy, MaxLogit, and KL), NAN shows significant improvement on both AUROC and FPR95. Moreover, NAN can be instantly applied to the contrastive models without a classification layer and label supervision.

Distance-based scores (Mahalanobis, SSD, and KNN) outperform NAN on the far-OOD dataset Texture. This is because NAN inherently is a classifier confidence, which can exhibit overconfidence when dealing with far OOD instances. On average, however, NAN is more robust and produces a significant reduction on the FPR95 metric (11-26%) without any hyperparameter. Rather than competing with the state-of-the-art distanced-based detectors, we show NAN can be integrated with them easily for further improvement.

6.2. Evaluation on NAN compatibility

We examine whether NAN can be integrated with existing OOD scores. To this end, we consider the state-of-the-art perturbation method ReAct [41] and the label-free distance-based scores SSD and KNN. NAN is combined with SSD and KNN by simple score division as follows: given a distance function to the ID bank set or prototypes in the form of $d(\mathbf{x}, X_{bank})$, we re-calibrate the distance by $d(\mathbf{x}, X_{bank}) / \|\mathbf{a}^{(L)}\|_{\text{NAN}}$ where $\mathbf{a}^{(L)}$ is the last hidden layer feature of the test input \mathbf{x} . Table 2 shows that the combination improves both metrics in all cases, demonstrating the compatibility of NAN.

	hyper.-free	label-free	bank-free	iNaturalist		SUN		Places		Texture		Average		ID ACC \uparrow
				AUROC \uparrow	FPR95 \downarrow									
With Supervised Labels of ID:														
MSP	✓		✓	93.78	29.74	84.56	59.54	84.28	60.94	84.90	50.02	86.88	50.06	78.73
Energy	✓		✓	96.17	20.98	88.91	47.05	87.70	51.15	88.90	39.31	90.42	39.62	78.73
MaxLogit	✓		✓	95.99	22.06	88.43	50.90	87.37	53.78	88.42	42.25	90.05	42.25	78.73
KL	✓		✓	96.17	20.98	88.91	47.06	87.70	51.15	88.90	39.31	90.42	39.63	78.73
Mahalanobis	✓		✓	94.79	35.04	86.55	64.99	83.92	70.31	95.52	15.02	90.20	46.34	78.73
ViM			✓	95.54	27.75	89.85	48.12	87.05	57.82	95.18	14.47	91.91	37.04	78.73
SSD		✓	✓	94.08	37.77	88.06	58.38	84.70	63.89	96.96	11.63	90.95	42.92	78.73
KNN		✓		94.15	38.25	87.75	58.19	84.93	61.80	94.24	19.29	90.27	44.38	78.73
NAN (ours)	✓		✓	96.94	15.86	92.77	29.81	91.46	37.21	88.09	43.46	92.32	31.59	78.73
Without Supervised Labels of ID (detectors based on supervised labels are not available):														
SSD		✓	✓	60.34	93.87	80.89	78.41	77.23	81.26	90.19	33.53	77.16	71.77	71.10
KNN		✓		84.53	78.71	82.26	76.06	77.50	80.65	91.99	24.61	84.07	65.01	71.10
NAN (ours)	✓	✓	✓	92.90	36.09	86.76	56.27	83.22	65.08	87.57	46.86	87.61	51.08	71.10

Table 1: Results on ImageNet-1k with ResNet-50. ‘hyper.-free’ indicates that the detection score does not require a hyperparameter.

	AUROC \uparrow	FPR95 \downarrow
NAN	92.32	31.59
NAN + KNN [42]	92.99	29.26
NAN + SSD [40]	93.42	27.51
NAN + ReAct [41]	93.91	29.23
NAN + ReAct [41] + KNN [42]	94.37	24.94
NAN + ReAct [41] + SSD [40]	94.61	24.57

Table 2: Compatibility of NAN to existing detectors. The ID is ImageNet-1k. The value is averaged over all test OOD datasets.

6.3. Evaluation on standard benchmark

We evaluate NAN on the standard CIFAR-10 benchmark that consists of low-resolution images.

Setup. We utilize a ResNet-18 trained on CIFAR10. The model is trained by either of the two standard training schemes: cross-entropy minimization with supervised labels and self-supervised learning (MoCo-v2) without the supervised labels. Following the popular benchmark, we choose the following datasets as OOD test datasets: LSUN-fix [43], ImageNet-fix [43], CIFAR100 [24], SVHN [34], and Places [56]. All images are of size 32×32 .

Evaluation results. Table 3 shows that the proposed score NAN is comparable to state-of-the-art scores specifically designed for OOD detection. We highlight that only NAN is a hyperparameter-free approach among the top-performing methods. The label-free distance-based scores KNN and SSD exhibit robustness, but their results are attained by carefully fine-tuning their method-specific hyperparameters. Despite not utilizing any hyperparameters, NAN exhibits comparable performance to the label-free state-of-the-art detectors (SSD and KNN) in terms of AUROC and FPR95 metrics on average. CSI also shows marginal superiority in two cases out of eight, but CSI requires complicated training with image rotation prediction, and its inference must be combined with KNN in an intricate manner. In contrast, NAN is simple and can be easily integrated to KNN. Combined with the distance-based scores SSD and KNN, NAN exhibits a consistent performance boost and outper-

forms all reported detectors.

6.4. Evaluation on one-class classification

As NAN requires neither classifier nor supervised labels, it can be applied to one-class classification (OCC). To assess the OCC performance, we evaluate the standard one-class benchmark of CIFAR-10/100. A class randomly chosen in CIFAR-10 is regarded as the ID data, and the rest of the 9 classes in CIFAR-10 constitute OOD instances. We conduct a similar experimental procedure on CIFAR-100 superclasses. For a fair comparison, we compare with one-class classification baselines that do not utilize extra training data and pretrained weights attained from large-scale data. For evaluation, we apply NAN on the MoCo-v2 model that is trained on the one-class data from scratch.

Table 4 indicates that NAN is comparable to the state-of-the-art one-class classifier CSI without any complicated training and hyperparameter tuning. Combined with the distance-based detectors, NAN performs equally well and improves the distance-based detectors on both CIFAR-10/100 data sets.

6.5. Ablation study

Ablation on the NAN Score The primary innovation of NAN is the inclusion of a sparsity term (*i.e.*, the denominator of NAN), which accounts for the hidden layer neurons’ tendency to deactivate. We analyze the impact of this component by ablating it. Table 5 shows the effectiveness of the sparsity term in both large-scale and small-scale settings. In the large-scale setting (ImageNet-1k), OOD is mostly differentiated from ID by the deactivation tendency of hidden layer neurons. In the case of the small-scale CIFAR-10 dataset, capturing both deactivation and activation tendencies is crucial for enhancing the OOD detection performance. In general, the inclusion of the sparsity term to capture the deactivation tendency enhances the robustness of the OOD detection score.

Ablation on the Architectural Component: the

OOD	LSUN-fix		ImageNet-fix		CIFAR100		SVHN		Places		Average		ID ACC \uparrow
	AUROC \uparrow	FPR95 \downarrow											
With supervised labels of ID													
ODIN* [26]	-	-	-	-	-	-	88.3	60.4	90.6	45.5	-	-	-
CSI* [43]	92.1	-	92.4	-	90.5	-	96.5	-	-	-	-	-	-
MSP	90.3	59.1	89.7	61.3	88.0	64.1	96.9	19.8	88.5	61.7	90.7	53.2	94.5
Energy	86.8	50.9	84.7	55.1	81.6	59.6	93.9	22.1	86.7	48.4	86.7	47.2	94.5
MaxLogit	86.8	51.7	84.7	56.0	81.6	60.1	94.1	22.0	86.6	49.8	86.8	47.9	94.5
KL	88.8	50.3	89.4	50.0	87.2	55.1	98.8	6.6	88.0	49.2	90.4	42.2	94.5
Mahalanobis	92.5	38.3	90.6	47.3	88.0	54.8	99.0	5.9	90.9	41.0	92.2	37.5	94.5
ViM	92.8	41.0	91.3	43.7	87.3	52.5	95.0	22.5	94.1	28.2	92.1	37.6	94.5
KNN	96.0	25.7	95.1	31.4	92.2	44.2	99.8	1.1	94.3	32.4	95.5	27.0	94.5
SSD	96.5	20.2	94.2	35.0	88.8	51.4	99.9	0.4	92.2	42.3	94.3	29.9	94.5
NAN (ours)	94.7	36.6	94.5	34.4	91.7	44.8	99.7	1.3	94.2	33.3	95.0	30.1	94.5
NAN + KNN	96.0	26.7	95.5	29.0	92.7	40.9	99.9	0.6	94.9	28.2	95.8	25.1	94.5
NAN + SSD	96.7	19.9	95.6	27.6	91.8	43.6	99.9	0.3	94.6	30.3	95.7	24.3	94.5
Without supervised labels of ID													
RotNet* [18]	81.6	-	86.7	-	82.3	-	97.8	-	-	-	-	-	-
GOAD* [1]	78.8	-	83.3	-	77.2	-	96.3	-	-	-	-	-	-
CSI* [43]	90.3	-	93.3	-	89.2	-	99.8	-	-	-	-	-	-
KNN	95.0	30.5	93.7	36.7	89.7	50.3	99.4	3.0	88.6	58.2	93.3	35.7	90.7
SSD	94.1	30.2	90.8	47.4	85.9	57.6	98.5	8.3	88.8	51.9	91.6	39.1	90.7
NAN (ours)	94.9	28.8	93.7	36.1	88.6	52.4	96.1	22.0	89.3	51.5	92.5	38.1	90.7
NAN + KNN	95.8	24.6	94.8	32.6	90.1	49.4	98.4	8.8	90.5	50.5	93.9	33.2	90.7
NAN + SSD	96.0	21.3	94.5	33.6	89.4	49.7	98.5	8.3	91.2	45.6	93.9	31.7	90.7

Table 3: **Results on CIFAR-10** with ResNet-18. * indicates the values are taken from the references.

	CIFAR10	CIFAR100
Without bank set:		
OC-SVM* [39]	58.8	63.1
Deep-SVDD* [37]	64.8	-
AnoGAN* [38]	61.8	-
OCGAN* [36]	65.7	-
Geom* [12]	86.0	78.7
GOAD* [2]	88.2	-
NAN (ours)	93.7	88.2
With bank set:		
CSI*	94.3	-
SSD	91.1	85.7
SSD + NAN (ours)	94.3(+3.2)	88.7(+2.0)
KNN	92.1	87.1
KNN + NAN (ours)	94.3(+2.2)	88.3(+1.0)

Table 4: The average one-class classification (OCC) performance in AUROC. * indicates the values are taken from the references.

	ImageNet-1k		CIFAR-10	
	AUROC \uparrow	FPR95 \downarrow	AUROC \uparrow	FPR95 \downarrow
NAN w/o sparsity term	57.99	95.22	92.40	43.00
NAN	92.32	31.59	94.90	30.10

Table 5: The ablation study examines the effect of NAN’s sparsity term, which accounts for the hidden layer neurons’ deactivation tendency. The ID is either ImageNet-1k or CIFAR-10. The value is averaged over all corresponding test OOD datasets.

Last Hidden Layer Dimension Although NAN is a hyperparameter-free OOD score, its effectiveness is still influenced by the network architecture, much like other detection scores. Specifically, the performance of NAN may primarily depend on the dimension d_L of the last hidden layer $\mathbf{a}^{(L)}$. To assess the impact of this dimension on the perfor-

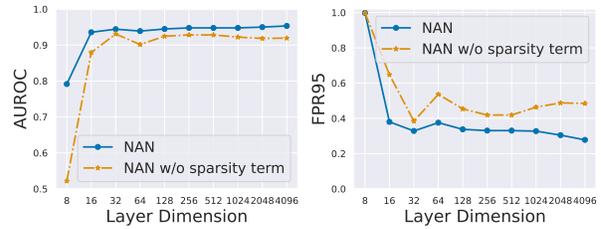


Figure 7: The ablation study of NAN with respect to the dimension d_L of the last hidden layer $\mathbf{a}^{(L)}$. The ID data is CIFAR-10. The reported metric numbers are values averaged over test OOD datasets.

formance of NAN, we evaluate NAN on multiple ResNet-18 models with different dimensions d_L . We train the models on CIFAR-10 using supervised cross-entropy loss and evaluate them on various OOD datasets, including LSUN-fix, ImageNet-fix, CIFAR-100, and SVHN. We report the average performance over all test OOD datasets.

We hypothesize that a wider hidden layer would better capture the deactivation tendency of neurons, and hence improve the performance. Fig. 7 evidences the hypothesis; increasing the dimension of the last hidden layer tends to improve the performance of NAN. Particularly on the FPR95 metric, the improvement is not marginal. Moreover, the performance is fairly robust unless the layer dimension is unreasonably small. Interestingly, the comparison between NAN and the standard l_1 -norm score without the sparsity term unveils an intriguing finding; NAN’s ability to capture the deactivation tendency makes the score more robust to changes in the layer dimension d_L . This result suggests that

measuring the deactivation tendency is critical for effective OOD detection.

Additional ablations and limitation. Further ablation on architectural components and the limitation of NAN are given in Sec. D and E, respectively.

7. Conclusion

We have conducted a thorough investigation of the feature norm to gain insights into its underlying mechanism for OOD detection. Specifically, we have demonstrated that the feature norm’s ability to detect OOD stems from its function as classifier confidence. Additionally, we have established that the feature norm can detect OOD using any discriminative model, making it independent of class label type. Through our formulation of the feature norm as a hidden classifier, we have identified that the conventional feature norm neglects neurons that tend to deactivate, leading to the potential misidentification of ID samples. To address this limitation, we have proposed a novel negative-aware norm NAN that captures both the activation and deactivation tendencies of hidden layer neurons. Our empirical results have demonstrated the effectiveness of NAN across diverse OOD detection benchmarks.

Acknowledgments This work was supported by the Materials/Parts Technology Development Program grant funded by the Korea government (MOTIE) (No. 1415187441) and the National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIP) (No. NRF-2022R1A2C1010710).

References

- [1] Liron Bergman and Yedid Hoshen. Classification-based anomaly detection for general data. *arXiv preprint arXiv:2005.02359*, 2020. 8
- [2] Liron Bergman and Yedid Hoshen. Classification-based anomaly detection for general data. *arXiv preprint arXiv:2005.02359*, 2020. 8
- [3] Di Chen, Shanshan Zhang, Jian Yang, and Bernt Schiele. Norm-aware embedding for efficient person search. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12615–12624, 2020. 1, 3
- [4] Di Chen, Shanshan Zhang, Jian Yang, and Bernt Schiele. Norm-aware embedding for efficient person search and tracking. *International Journal of Computer Vision*, 129(11):3154–3168, 2021. 3
- [5] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020. 6, 12, 16
- [6] Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3606–3613, 2014. 4, 6
- [7] Akshay Raj Dhamija, Manuel Günther, and Terrance Boult. Reducing network agnostophobia. *Advances in Neural Information Processing Systems*, 31, 2018. 1, 2
- [8] Thomas G Dietterich and Alex Guyer. The familiarity hypothesis: Explaining the behavior of deep open set methods. *Pattern Recognition*, 132:108931, 2022. 1, 21
- [9] Zhen Fang, Yixuan Li, Jie Lu, Jiahua Dong, Bo Han, and Feng Liu. Is out-of-distribution detection learnable? *arXiv preprint arXiv:2210.14707*, 2022. 1, 3, 4, 5, 21
- [10] Zhen Fang, Jie Lu, Anjin Liu, Feng Liu, and Guangquan Zhang. Learning bounds for open-set learning. In *International Conference on Machine Learning*, pages 3122–3132. PMLR, 2021. 1, 2, 4
- [11] Kunihiko Fukushima. Cognitron: A self-organizing multi-layered neural network. *Biological cybernetics*, 20(3):121–136, 1975. 3
- [12] Izhak Golan and Ran El-Yaniv. Deep anomaly detection using geometric transformations. *Advances in neural information processing systems*, 31, 2018. 8
- [13] Will Grathwohl, Kuan-Chieh Wang, Jörn-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi, and Kevin Swersky. Your classifier is secretly an energy based model and you should treat it like one. *arXiv preprint arXiv:1912.03263*, 2019. 2
- [14] Boris Hanin and David Rolnick. Deep relu networks have surprisingly few activation patterns. *Advances in neural information processing systems*, 32, 2019. 4
- [15] Dan Hendrycks, Steven Basart, Mantas Mazeika, Mohammadreza Mostajabi, Jacob Steinhardt, and Dawn Song. Scaling out-of-distribution detection for real-world settings. *arXiv preprint arXiv:1911.11132*, 2019. 2
- [16] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016. 1, 2
- [17] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016. 3
- [18] Dan Hendrycks, Mantas Mazeika, Saurav Kadavath, and Dawn Song. Using self-supervised learning can improve model robustness and uncertainty. *Advances in neural information processing systems*, 32, 2019. 8
- [19] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15262–15271, 2021. 1
- [20] Yen-Chang Hsu, Yilin Shen, Hongxia Jin, and Zsolt Kira. Generalized odin: Detecting out-of-distribution image without learning from out-of-distribution data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10951–10960, 2020. 2
- [21] Rui Huang, Andrew Geng, and Yixuan Li. On the importance of gradients for detecting distributional shifts in the wild. *Advances in Neural Information Processing Systems*, 34:677–689, 2021. 3
- [22] Rui Huang and Yixuan Li. Mos: Towards scaling out-of-distribution detection for large semantic space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8710–8719, 2021. 6

- [23] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschiot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in Neural Information Processing Systems*, 33:18661–18673, 2020. 6, 18
- [24] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 4, 7
- [25] Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. *Advances in neural information processing systems*, 31, 2018. 2
- [26] Shiyu Liang, Yixuan Li, and Rayadurgam Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. *arXiv preprint arXiv:1706.02690*, 2017. 2, 8
- [27] Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection. *Advances in Neural Information Processing Systems*, 33:21464–21475, 2020. 2
- [28] Qiang Meng, Shichao Zhao, Zhida Huang, and Feng Zhou. Magface: A universal representation for face recognition and quality assessment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14225–14234, 2021. 1, 3
- [29] Yifei Ming, Yiyou Sun, Ousmane Dia, and Yixuan Li. How to exploit hyperspherical embeddings for out-of-distribution detection? *arXiv preprint arXiv:2203.04450*, 2022. 2, 20
- [30] Yifei Ming, Hang Yin, and Yixuan Li. On the impact of spurious correlation for out-of-distribution detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 10051–10059, 2022. 1
- [31] Peyman Morteza and Yixuan Li. Provable guarantees for understanding out-of-distribution detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 8, 2022. 1
- [32] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *icml*, 2010. 3
- [33] Eric Nalisnick, Akihiro Matsukawa, Yee Whye Teh, Dilan Gorur, and Balaji Lakshminarayanan. Do deep generative models know what they don’t know? *arXiv preprint arXiv:1810.09136*, 2018. 1
- [34] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bisaccho, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011. 4, 7
- [35] Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 427–436, 2015. 1
- [36] Pramuditha Perera, Ramesh Nallapati, and Bing Xiang. Ocgan: One-class novelty detection using gans with constrained latent representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2898–2906, 2019. 8
- [37] Lukas Ruff, Robert Vandermeulen, Nico Goernitz, Lucas Deecke, Shoaib Ahmed Siddiqui, Alexander Binder, Emmanuel Müller, and Marius Kloft. Deep one-class classification. In *International conference on machine learning*, pages 4393–4402. PMLR, 2018. 8
- [38] Thomas Schlegl, Philipp Seeböck, Sebastian M Waldstein, Ursula Schmidt-Erfurth, and Georg Langs. Unsupervised anomaly detection with generative adversarial networks to guide marker discovery. In *Information Processing in Medical Imaging: 25th International Conference, IPMI 2017, Boone, NC, USA, June 25-30, 2017, Proceedings*, pages 146–157. Springer, 2017. 8
- [39] Bernhard Schölkopf, Robert C Williamson, Alex Smola, John Shawe-Taylor, and John Platt. Support vector method for novelty detection. *Advances in neural information processing systems*, 12, 1999. 8
- [40] Vikash Schwag, Mung Chiang, and Prateek Mittal. Ssd: A unified framework for self-supervised outlier detection. *arXiv preprint arXiv:2103.12051*, 2021. 2, 7
- [41] Yiyou Sun, Chuan Guo, and Yixuan Li. React: Out-of-distribution detection with rectified activations. *Advances in Neural Information Processing Systems*, 34:144–157, 2021. 2, 6, 7
- [42] Yiyou Sun, Yifei Ming, Xiaojin Zhu, and Yixuan Li. Out-of-distribution detection with deep nearest neighbors. *arXiv preprint arXiv:2204.06507*, 2022. 2, 7, 18, 20
- [43] Jihoon Tack, Sangwoo Mo, Jongheon Jeong, and Jinwoo Shin. Csi: Novelty detection via contrastive learning on distributionally shifted instances. *Advances in neural information processing systems*, 33:11839–11852, 2020. 2, 3, 7, 8, 18
- [44] Grant Van Horn, Oisín Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8769–8778, 2018. 6
- [45] Sagar Vaze, Kai Han, Andrea Vedaldi, and Andrew Zisserman. Open-set recognition: A good closed-set classifier is all you need. *arXiv preprint arXiv:2110.06207*, 2021. 1, 2, 3
- [46] Haoqi Wang, Zhizhong Li, Litong Feng, and Wayne Zhang. Vim: Out-of-distribution with virtual-logit matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4921–4930, 2022. 2
- [47] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *2010 IEEE computer society conference on computer vision and pattern recognition*, pages 3485–3492. IEEE, 2010. 6
- [48] Bing Xu, Naiyan Wang, Tianqi Chen, and Mu Li. Empirical evaluation of rectified activations in convolutional network. *arXiv preprint arXiv:1505.00853*, 2015. 3
- [49] Pingmei Xu, Krista A Ehinger, Yinda Zhang, Adam Finkelstein, Sanjeev R Kulkarni, and Jianxiong Xiao. Turkergaze: Crowdsourcing saliency with webcam based eye tracking. *arXiv preprint arXiv:1504.06755*, 2015. 4
- [50] Ruijia Xu, Guanbin Li, Jihan Yang, and Liang Lin. Larger norm more transferable: An adaptive feature norm approach for unsupervised domain adaptation. In *Proceedings of the*

- IEEE/CVF International Conference on Computer Vision*, pages 1426–1435, 2019. 3
- [51] Jingkang Yang, Pengyun Wang, Dejian Zou, Zitang Zhou, Kunyuan Ding, Wenxuan Peng, Haoqi Wang, Guangyao Chen, Bo Li, Yiyun Sun, et al. Openood: Benchmarking generalized out-of-distribution detection. *Advances in Neural Information Processing Systems*, 35:32598–32611, 2022. 1, 21
- [52] Jingkang Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu. Generalized out-of-distribution detection: A survey. *arXiv preprint arXiv:2110.11334*, 2021. 1
- [53] Chang Yu, Xiangyu Zhu, Zhen Lei, and Stan Z Li. Out-of-distribution detection for reliable face recognition. *IEEE Signal Processing Letters*, 27:710–714, 2020. 1, 3
- [54] Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015. 4
- [55] Yuhui Yuan, Kuiyuan Yang, and Chao Zhang. Feature incay for representation regularization. *arXiv preprint arXiv:1705.10284*, 2017. 3
- [56] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE transactions on pattern analysis and machine intelligence*, 40(6):1452–1464, 2017. 4, 6, 7

Supplementary Materials

A. Supplementary to the Analysis of Hidden Classifier

A.1. Proofs for the properties of hidden classifier

Notation (detailed) Each hidden layer feature $\mathbf{a}^{(l)}$ is defined by consecutive computation of the post-activated feature vector

$$\mathbf{a}^{(l)} = \sigma(\mathbf{W}^{(l)T} \mathbf{a}^{(l-1)}) \quad (11)$$

from the input layer $l = 0$ to the last hidden layer $l = L$. The pre-activated features satisfy $\mathbf{a}^{(l)} = \sigma(\mathbf{z}^{(l)})$, where the activation function σ is a rectifier (e.g. ReLU, GeLU, Leaky ReLU). The penultimate embedding is $g(\mathbf{x}) = \mathbf{U}^T \mathbf{a}^{(L)}$, which computes the network classification logit $\psi(\mathbf{x}) \in \mathbb{R}^K$ by

$$\psi(\mathbf{x}) = \mathbf{W}^T g(\mathbf{x}). \quad (12)$$

\mathbf{W} is the weight matrix for the classification layer. For notation simplicity, let $\mathbf{W}^{(L+1)} := \mathbf{U}\mathbf{W}$ such that $\psi(\mathbf{x}) = \mathbf{W}^{(L+1)} \mathbf{a}^{(L)}$. The sign function $\text{sign}(\cdot)$, on the other hand, that binarizes a scalar to either 1 or -1 is applied point-wise.

Note For the embedding computation, \mathbf{U} is a fixed identity matrix in supervised models, while \mathbf{U} serves as a learnable parameters for self-supervised models with projection head [5].

Proposition 1. *The final logit is represented by*

$$\psi(\mathbf{x}) = \mathbf{C}^{(l)} \mathbf{a}^{(l)} \quad (13)$$

for each hidden layer l , where

$$\mathbf{C}^{(l)} = \left(\prod_{k=0}^{L-l-1} \mathbf{W}^{(L+1-k)T} \mathbf{D}^{(L-k)} \right) \mathbf{W}^{(l+1)T} \quad (14)$$

with $\mathbf{D}^{(l)} = \text{diag}(\frac{\sigma(z_1)}{z_1}, \dots, \frac{\sigma(z_{d_l})}{z_{d_l}})$ with the convention $\frac{\sigma}{0} = 0$. $\mathbf{C}^{(l)} = \mathbf{C}^{(l)}(\mathbf{x}) \in \mathbb{R}^{K \times d_l}$ depends on \mathbf{x} .

Proof. Observe inductively that

$$\psi(\mathbf{x}) = \mathbf{W}^{(L+1)T} \mathbf{a}^{(L)} \quad (15)$$

$$= \mathbf{W}^{(L+1)T} \mathbf{D}^{(L)} \mathbf{z}^{(L)} \quad (16)$$

$$= \mathbf{W}^{(L+1)T} \mathbf{D}^{(L)} \mathbf{W}^{(L)T} \mathbf{a}^{(L-1)} \quad (17)$$

$$= \mathbf{W}^{(L+1)T} \mathbf{D}^{(L)} \mathbf{W}^{(L)T} \mathbf{D}^{(L-1)} \mathbf{z}^{(L-1)} \quad (18)$$

$$= \dots, \quad (19)$$

obtaining

$$\psi(\mathbf{x}) = \left(\prod_{k=0}^{L-l-1} \mathbf{W}^{(L+1-k)T} \mathbf{D}^{(L-k)} \right) \mathbf{W}^{(l+1)T} \mathbf{a}^{(l)} \quad (20)$$

□

Remark. We note that both $\mathbf{D}^{(l)} = \mathbf{D}^{(l)}(\mathbf{x})$ and $\mathbf{C}^{(l)} = \mathbf{C}^{(l)}(\mathbf{x})$ depend on \mathbf{x} as they depend on $\mathbf{a}^{(l)}$. Also, note that the dimension of $\mathbf{C}^{(l)}$ is $K \times d_l$.

Recall that $\mathbf{C}^{(l)} = [\mathbf{c}_1^{(l)}, \dots, \mathbf{c}_K^{(l)}]^T$.

Proposition 2. *Let (\mathbf{x}, y) be arbitrary labeled ID sample. Suppose that $\psi_y(\mathbf{x})$ is maximized in a manner to reduce the angle between $\mathbf{c}_y^{(l)}$ and $\mathbf{a}^{(l)}$ sufficiently that $\text{sign}(\mathbf{c}_y^{(l)}) = \text{sign}(\mathbf{a}^{(l)})$. Suppose that $\psi_k(\mathbf{x})$ is minimized in a manner to increase the angle between $\mathbf{c}_k^{(l)}$ and $\mathbf{a}^{(l)}$ sufficiently that $\angle(\text{sign}(\mathbf{c}_k^{(l)}), \mathbf{a}^{(l)}) > \pi/2$. Then, $\bar{\psi}^{(l)}$ becomes a discriminative classifier with $\bar{\psi}_y^{(l)}(\mathbf{x}) > \bar{\psi}_k^{(l)}(\mathbf{x})$.*

Proof. For notational simplicity, ignore the superscript index l , and let $\mathbf{a} = \mathbf{a}^{(l)}$, $\mathbf{b}_k = \mathbf{b}_k^{(l)}$, $\mathbf{c}_k = \mathbf{c}_k^{(l)}$, and $\bar{\psi} = \bar{\psi}^{(l)}$. First, observe $\mathbf{b}_y = \text{sign}(\mathbf{c}_y) = \text{sign}(\mathbf{a})$ implies $0 \leq \angle(\mathbf{b}_y, \mathbf{a}) < \pi/2$. Therefore,

$$\bar{\psi}_y(\mathbf{x}) = \mathbf{b}_y \cdot \mathbf{a} = \|\mathbf{b}_y\|_2 \|\mathbf{a}\|_2 \cos(\angle(\mathbf{b}_y, \mathbf{a})) > 0. \quad (21)$$

On the other hand, $\angle(\text{sign}(\mathbf{c}_k), \mathbf{a}) > \pi/2$ means $\pi \geq \angle(\mathbf{b}_k, \mathbf{a}) > \pi/2$ by the definition of \mathbf{b}_k for $k \neq y$. Therefore,

$$\bar{\psi}_k(\mathbf{x}) = \mathbf{b}_k \cdot \mathbf{a} = \|\mathbf{b}_k\|_2 \|\mathbf{a}\|_2 \cos(\angle(\mathbf{b}_k, \mathbf{a})) < 0. \quad (22)$$

Since (\mathbf{x}, y) was arbitrary, we have proved the desired. \square

The main message of Prop. 2 is that the discriminative optimization of the original classifier should be powerful enough to optimize the *angle* between the hidden layer feature and the binary weight. Then, in this case, the hidden classifier becomes discriminative.

Theorem 3. Under the sufficient condition of Prop. 2, for any labeled ID sample (\mathbf{x}, y) ,

$$\|\mathbf{a}^{(l)}\|_1 \text{ converges to } \bar{\psi}_y^{(l)}(\mathbf{x}) = \max_k \bar{\psi}_k^{(l)}(\mathbf{x}) \quad (23)$$

in which case $\text{sign}(\mathbf{a}^{(l)}) = \mathbf{b}_y^{(l)}$. In general, for any k and for any sample \mathbf{x} (either ID or OOD),

$$0 \leq \|\mathbf{a}^{(l)}\|_1 - \bar{\psi}_k^{(l)}(\mathbf{x}) \leq \|\mathbf{a}^{(l)}\|_\infty \|\text{sign}(\mathbf{a}^{(l)}) - \mathbf{b}_k^{(l)}\|_1. \quad (24)$$

Proof. For notational simplicity, ignore the superscript index l , and let $\mathbf{a} = \mathbf{a}^{(l)}$, $\mathbf{b}_k = \mathbf{b}_k^{(l)}$, $\mathbf{c}_k = \mathbf{c}_k^{(l)}$, and $\bar{\psi} = \bar{\psi}^{(l)}$. First, observe that

$$\|\mathbf{a}\|_1 = \sum_i |a_i| \geq \sum_i b_{ki} a_i = \mathbf{b}_k \cdot \mathbf{a} = \bar{\psi}_k(\mathbf{x}) \quad (25)$$

where $\mathbf{b}_k = (b_{k1}, \dots, b_{kd_l}) \in \{-1, 1\}^{d_l}$. This proves that $\|\mathbf{a}\|_1 \geq \bar{\psi}_k(\mathbf{x})$ for all k .

Now, observe that $|a_i| = \text{sign}(a_i) a_i$. Therefore,

$$\|\mathbf{a}\|_1 - \bar{\psi}_k(\mathbf{x}) = \sum_i (\text{sign}(a_i) - b_{ki}) a_i \leq \sum_i |\text{sign}(a_i) - b_{ki}| |a_i| \leq \|\mathbf{a}\|_\infty \|\text{sign}(\mathbf{a}) - \mathbf{b}_k\|_1, \quad (26)$$

proving a general upper bound of the difference between the hidden classifier output and the feature norm.

Now, under the sufficient condition of Prop. 2, the binary weight becomes the activation pattern by the assumption; $\text{sign}(\mathbf{a}) = \mathbf{b}_y$. Therefore, in this case,

$$0 \leq \|\mathbf{a}\|_1 - \bar{\psi}_y(\mathbf{x}) \leq \|\mathbf{a}\|_\infty \cdot 0 = 0, \quad (27)$$

proving the desired. \square

Corollary 4. If $\max_k \bar{\psi}_k^{(l)}(\mathbf{x}_{ood})$ is sufficiently small such that

$$\max_k \bar{\psi}_k^{(l)}(\mathbf{x}_{ood}) + \delta < \max_k \bar{\psi}_k^{(l)}(\mathbf{x}_{ind}) \quad (28)$$

for all ID samples \mathbf{x}_{ind} where

$$\delta \geq \|\mathbf{a}^{(l)}(\mathbf{x}_{ood})\|_\infty \cdot \|\text{sign}(\mathbf{a}^{(l)}(\mathbf{x}_{ood})) - \mathbf{b}_{k_0}^{(l)}\|_1 \quad (29)$$

and $k_0 = \arg \max_k \bar{\psi}_k^{(l)}(\mathbf{x}_{ood})$, then

$$\|\mathbf{a}^{(l)}(\mathbf{x}_{ood})\|_1 < \|\mathbf{a}^{(l)}(\mathbf{x}_{ind})\|_1 \quad (30)$$

for all ID samples \mathbf{x}_{ind} .

Proof. By Thm. 3,

$$\|\mathbf{a}^{(l)}(\mathbf{x}_{ood})\|_1 \leq \bar{\psi}_{k_0}^{(l)}(\mathbf{x}_{ood}) + \|\mathbf{a}^{(l)}(\mathbf{x}_{ood})\|_\infty \|\text{sign}(\mathbf{a}^{(l)}(\mathbf{x}_{ood})) - \mathbf{b}_{k_0}^{(l)}\|_1 < \max_k \bar{\psi}_k^{(l)}(\mathbf{x}_{ind}) \leq \|\mathbf{a}^{(l)}(\mathbf{x}_{ind})\|_1. \quad (31)$$

\square

A.2. Additional Theoretical Consideration

We present additional results of the theoretical analysis on the hidden classifier.

A.2.1 Relation to General l_p -norms

We have proved that l_1 -norm can differentiate OOD from ID. This capability of l_1 -norm extends to the general l_p -norm by Holder's inequality.

Theorem 5 (Holder's inequality). *For $0 < p \leq q < \infty$ and $\mathbf{x} \in \mathbb{R}^d$,*

$$\|\mathbf{x}\|_q \leq \|\mathbf{x}\|_p \leq d^{1/p-1/q} \|\mathbf{x}\|_q. \quad (32)$$

Thus, for an activation vector $\mathbf{a}^{(l)} \in \mathbb{R}^{d_l}$ and for $p > 1$, we have

$$\|\mathbf{a}^{(l)}\|_p \leq \|\mathbf{a}^{(l)}\|_1 \leq d_l^{-1/p} \|\mathbf{a}^{(l)}\|_p \quad (33)$$

Therefore, if $\|\mathbf{a}^{(l)}\|_1$ is large or small, then $\|\mathbf{a}^{(l)}\|_p$ is also large or small, respectively. Thus, different l_p -norms have similar mechanisms for OOD detection. Note, however, that different l_p -norms have different priors on the computation of units in the activation vector. Accordingly, the OOD detection performance will vary depending on which l_p -norm is used.

A.2.2 Extension to Pre-Activation Layer

Extending the framework in Sec. 3 to the pre-activation layer feature vector $\mathbf{z}^{(l)}$ is trivial, where the pre-activation layer feature is the vector satisfying $\mathbf{a}^{(l)} = \sigma(\mathbf{z}^{(l)})$ with the activation function σ . Here, we provide the properties of the pre-activation layer that correspond to the ones given in Sec. 3.

Proposition 6. *The final logit is represented by*

$$\psi(\mathbf{x}) = \widehat{\mathbf{C}}^{(l)} \mathbf{z}^{(l)} \quad (34)$$

for each hidden layer l , where

$$\widehat{\mathbf{C}}^{(l)} = \left(\prod_{k=0}^{L-l} \mathbf{W}^{(L+1-k)T} \mathbf{D}^{(L-k)} \right) \quad (35)$$

with $\mathbf{D}^{(l)} = \text{diag}(\frac{\sigma(z_{d_1})}{z_{d_1}}, \dots, \frac{\sigma(z_{d_l})}{z_{d_l}})$ with the convention $\frac{\sigma(z_{d_l})}{z_{d_l}} = 0$. $\widehat{\mathbf{C}}^{(l)} = \widehat{\mathbf{C}}^{(l)}(\mathbf{x}) \in \mathbb{R}^{K \times d_l}$ depends on \mathbf{x} .

Define a hidden classifier corresponding to $\mathbf{z}^{(l)}$ by

$$\widehat{\psi}(\mathbf{x}) := \text{sign}(\widehat{\mathbf{C}}^{(l)} \mathbf{z}^{(l)}) = \widehat{\mathbf{B}}^{(l)} \mathbf{z}^{(l)} \quad (36)$$

where $\widehat{\mathbf{C}}^{(l)} = [\widehat{\mathbf{c}}_1^{(l)}, \dots, \widehat{\mathbf{c}}_K^{(l)}]^T$ and $\widehat{\mathbf{B}}^{(l)} = [\widehat{\mathbf{b}}_1^{(l)}, \dots, \widehat{\mathbf{b}}_K^{(l)}]^T$.

Proposition 7. *Let (\mathbf{x}, y) be an arbitrary labeled sample. Suppose that $\psi_y(\mathbf{x})$ is maximized in a manner to reduce the angle between $\widehat{\mathbf{c}}_y^{(l)}$ and $\mathbf{z}^{(l)}$ sufficiently that $\text{sign}(\widehat{\mathbf{c}}_y^{(l)}) = \text{sign}(\mathbf{z}^{(l)})$. Suppose that $\psi_k(\mathbf{x})$ is minimized in a manner to increase the angle between $\widehat{\mathbf{c}}_k^{(l)}$ and $\mathbf{z}^{(l)}$ sufficiently that $\angle(\text{sign}(\widehat{\mathbf{c}}_k^{(l)}), \mathbf{z}^{(l)}) > \pi/2$. Then, $\widehat{\psi}^{(l)}$ becomes a discriminative classifier with $\widehat{\psi}_y^{(l)}(\mathbf{x}) > \widehat{\psi}_k^{(l)}(\mathbf{x})$.*

Theorem 8. *Under the sufficient condition of Prop. 7,*

$$\|\mathbf{z}^{(l)}\|_1 \text{ converges to } \widehat{\psi}_y^{(l)}(\mathbf{x}) = \max_k \widehat{\psi}_k^{(l)}(\mathbf{x}) \quad (37)$$

in which case $\text{sign}(\mathbf{z}^{(l)}) = \widehat{\mathbf{b}}_y^{(l)}$. In general, for any k

$$0 \leq \|\mathbf{z}^{(l)}\|_1 - \widehat{\psi}_k(\mathbf{x}) \leq \|\mathbf{z}^{(l)}\|_\infty \|\text{sign}(\mathbf{z}^{(l)}) - \widehat{\mathbf{b}}_k^{(l)}\|_1 \quad (38)$$

A.2.3 On Bias

In Sec. 3, we ignored the bias in the computation of features for simplicity. We can preserve the properties of features given in Sec. 3 while including the bias terms. To observe this, consider

$$\mathbf{a}^{(l)} = \sigma(\mathbf{W}^{(l)T} \mathbf{a}^{(l-1)} + \boldsymbol{\beta}^{(l)}) = \mathbf{D}^{(l)} \mathbf{W}^{(l)T} \mathbf{a}^{(l-1)} + \mathbf{D}^{(l)} \boldsymbol{\beta}^{(l)}. \quad (39)$$

Thus, if Ψ denotes the logit computed with bias, then

$$\Psi(x) = \mathbf{C}^{(l)} \mathbf{a}^{(l)} + \sum_{j=l}^L \widehat{\mathbf{C}}^{(j+1)} \boldsymbol{\beta}^{(j+1)} = \psi(x) + \boldsymbol{\Gamma} \quad (40)$$

with $\boldsymbol{\Gamma} = \boldsymbol{\Gamma}(l, \mathbf{x}) = \sum_{j=l}^L \widehat{\mathbf{C}}^{(j+1)} \boldsymbol{\beta}^{(j+1)}$ and the convention that $\widehat{\mathbf{C}}^{(L+1)} = \mathbf{I}$. Hence, if the discriminative learning of Ψ is not trivially achieved by the optimization of the bias term $\boldsymbol{\Gamma}$, and if the discriminative learning of ψ is thus sufficiently powerful, then the properties in Sec. 3 hold.

A.2.4 On Cosine Similarity Logit

We assumed that the classification logit is the output of the inner product in Sec. 3. Here, we show that changing the inner product logit by a (scaled) cosine similarity logit does not alter the major behavior of discriminative learning, and hence they are equivalent in our theoretical consideration. Thus, the theory developed in the inner-product logit also holds in the (scaled) cosine similarity logit.

To observe this, note that the scaled cosine similarity logit is defined as

$$\phi_k(\mathbf{x}) = \frac{1}{T} \frac{\mathbf{w}_k \cdot g(\mathbf{x})}{\|\mathbf{w}_k\|_2 \|g(\mathbf{x})\|_2} \quad (41)$$

where \mathbf{w}_k are class weight vectors (prototypes) of trainable parameters and $g(\mathbf{x}) = \mathbf{U}^T \mathbf{a}^{(L)}$ with a matrix \mathbf{U} of trainable parameters. T is the temperature that modifies the scale of similarity. Without loss of generality, we assume $T = 1$. Let $\psi_k(\mathbf{x}) = \mathbf{w}_k \cdot g(\mathbf{x})$ denote the inner-product logit that we originally used. Thus, we have

$$\phi_k(x) = \psi_k(x) (\|\mathbf{w}_k\|_2 \|g(\mathbf{x})\|_2)^{-1}. \quad (42)$$

During discriminative learning, the model maximizes

$$(-1)^{1_{y \neq k}} \phi_k(\mathbf{x}) = (-1)^{1_{y \neq k}} \psi_k(\mathbf{x}) (\|\mathbf{w}_k\|_2 \|g(\mathbf{x})\|_2)^{-1}. \quad (43)$$

Assuming $\psi_y(\mathbf{x}) = \mathbf{w}_y \cdot g(\mathbf{x}) > 0$ and $\psi_k(\mathbf{x}) = \mathbf{w}_k \cdot g(\mathbf{x}) < 0$, the above maximization is equivalent to minimizing its negative log

$$-\log((-1)^{1_{y \neq k}} \phi_k(\mathbf{x})) = -\log((-1)^{1_{y \neq k}} \psi_k(\mathbf{x})) + \log(\|\mathbf{w}_k\|_2 \|g(\mathbf{x})\|_2), \quad (44)$$

which can be considered as the constrained minimization of

$$-\log((-1)^{1_{y \neq k}} \psi_k(\mathbf{x})) \equiv -(-1)^{1_{y \neq k}} \psi_k(\mathbf{x}) \quad (45)$$

constraint to

$$\|\mathbf{w}_k\|_2 \|g(\mathbf{x})\|_2 \leq e^{\eta_0} = \eta \quad (46)$$

for some η . Thus, optimization of the cosine similarity logit is equivalent to the constrained optimization of the inner product logit.

Proposition 9. *The maximization*

$$\max_{\phi} (-1)^{1_{y \neq k}} \phi_k(\mathbf{x}) \quad (47)$$

is equivalent to

$$\begin{aligned} & \max_{\psi} (-1)^{1_{y \neq k}} \psi_k(\mathbf{x}) \\ & \text{subject to } \|\mathbf{w}_k\|_2 \|g(\mathbf{x})\|_2 \leq \eta \end{aligned} \quad (48)$$

for some $\eta > 0$ if $\psi_y > 0$ and $\psi_k < 0$.

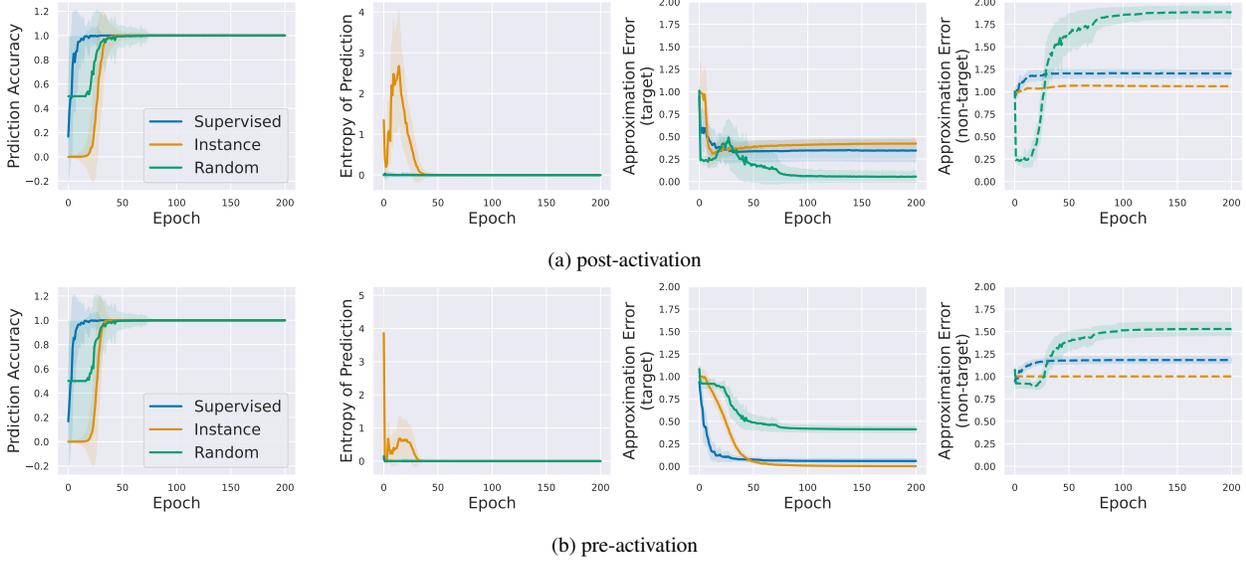


Figure 8: Results of hidden classifiers of ResNet-18 with different class labeling schemes on CIFAR-10. The approximation error on the target unit measures the normalized error $(\|\mathbf{a}\|_1 - \bar{\psi}_y(\mathbf{x})) / \|\mathbf{a}\|_1$, while the approximation error on the non-target unit is the average of $(\|\mathbf{a}\|_1 - \bar{\psi}_k(\mathbf{x})) / \|\mathbf{a}\|_1$ with respect to $k \neq y$. In the case of post-activation, the vector \mathbf{a} is $\mathbf{a} = \mathbf{a}^{(L)}$. In the case of pre-activation, the vector \mathbf{a} is $\mathbf{a} = \mathbf{z}^{(L)}$.

A.3. Supplementary to empirical validation of hidden classifier

Here, we provide a detailed description of the experiments conducted to validate the theoretical analysis presented in Sec. 3.

A.3.1 On MLP

Setup We train an MLP with 5 hidden layers. The hidden layer dimension is fixed to 512, and likewise for the embedding layer dimension. The embedding is normalized, and the cosine similarity logit is divided by a temperature of 0.1. The model is trained by AdamW for 200 epochs with batch size 256. The learning rate decays from 0.001 to 0 by the cosine scheduler. Other setups follow the default setting in PyTorch.

Results The results are given in Fig. 12, 13, and 14. They have similar trends that we expected and thus verify our theoretical claims.

A.3.2 On Convolutional Network

Setup The experiment setup is given as in Sec. B.

Results In the cases of both instance discrimination (I), supervised learning (S), and random binary label discrimination (R), the hidden classifier of the last hidden layer in ResNet-18 is trained to be discriminative (Fig. 8).

B. Supplementary to the Analysis of Feature Norm’s Class Agnosticity

Setup. We train a ResNet-18 on CIFAR-10. We add an MLP projection head as in MoCo-v2 [5]. The embedding is normalized, and the cosine similarity logit is divided by a temperature of 0.1. The model is trained for 200 epochs and batch size 256 with the SGD optimizer, cosine learning rate (0.06 to 0), and momentum 0.9. Each model is trained in a different manner based on a different class labeling scheme:

- **S:** The class labels y_i are supervised labels (e.g. plane, dog, cat, ...). No data augmentation is applied.

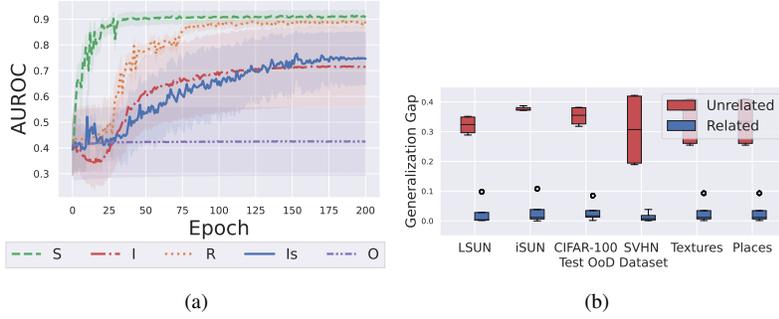


Figure 9: (a) The detection performance of NAN versus the learning epoch across different types of training schemes (b) The generalization gap of NAN based on the intra-class semantics.

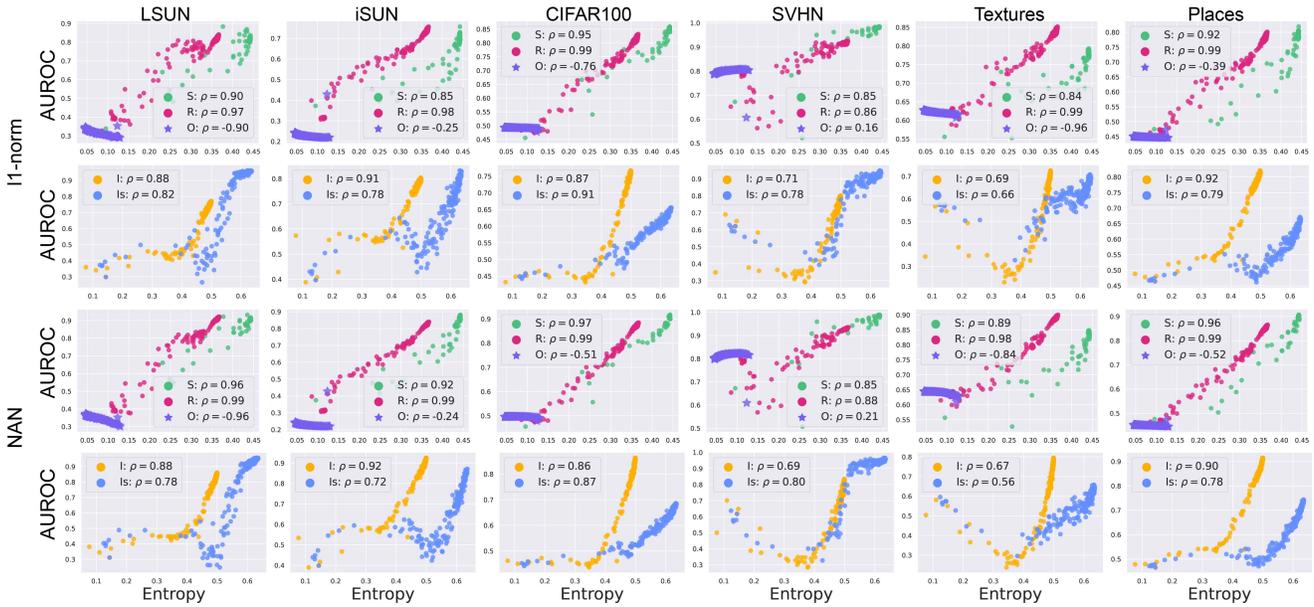


Figure 10: The graph of the detection performance versus the activation entropy. The performance is measured at every training epoch.

- **I**: The class labels y_i are instance labels $y_i = i$. No data augmentation is applied such that each instance class has only one intra-class sample.
- **Is**: The class labels y_i are instance labels $y_i = i$. Data augmentation is applied such that each instance class has multiple intra-class samples.
- **R**: The class labels y_i are labeled randomly by a binary number $y_i \in \{0, 1\}$.
- **O**: The class labels y_i are labeled with a single label $y_i = 0$ such that every sample is in the same class.

Other setups follow the default setting in PyTorch.

Full results on the impact of inter/intra-class learning The additional results on NAN is given in Fig. 9, which NAN exhibits the same trend of memorization and generalization as the conventional feature norm.

Full results on the relation to entropy The full results on the relation between the activation entropy and the detection performance is given in Fig. 10.

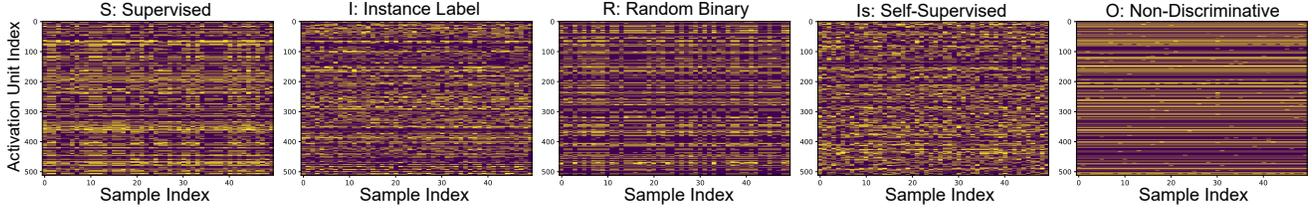


Figure 11: **Activation patterns** of randomly chosen 50 ID samples after training. Each column corresponds to the activation pattern $\text{sign}(\mathbf{a}^{(L)}) \in \{-1, 1\}^{512}$ of an ID sample. Discriminative training $\{S, R, I, Is\}$ results in *diverse* activation patterns, while the activation pattern *collapses* for the non-discriminative model O.

On the activation pattern If the model is trained in a non-discriminative manner with a single class, then the entropy of activation is diminished. In this case, the activation pattern collapses as shown in Fig. 11.

C. The Detailed Setup for the Experiments on NAN

C.1. Setup

Setup: ImageNet-1k For the supervised model trained by the cross entropy, we utilize the ResNet-50 backbone trained on ImageNet-1k. The model is provided by the PyTorch model zoo.

For the supervised model trained by the contrastive loss (thanks to the authors of [42]), we utilize the pretrained ResNet-50 model provided from the official GitHub page of KNN [42], which is trained on ImageNet-1k by the supervised contrastive loss [23] with the MLP projection head.

For the self-supervised contrastive model trained without the supervised labels of ID, thanks to the authors of MoCo-v2, we utilize the pretrained MoCo-v2 model provided from the official GitHub page of MoCo-v2 (the one with 71.1 accuracies on ImageNet-1k).

Setup: OOD CIFAR-10 For the evaluation results of OOD detection ‘with supervised labels of ID’ in Table 3, we train a cross-entropy model with supervised labels of CIFAR-10. The model has trained on CIFAR-10 over 800 epochs with the SGD optimizer and its momentum is 0.9. The learning rate decays to 0 from 0.03 by the cosine scheduler. The batch size is 512. The backbone is ResNet-18, accompanied by an MLP projection head on top of the encoder as in MoCo-v2. The embedding is normalized, and the cosine similarity logit is divided by the temperature 0.1.

For the evaluation results of OOD detection ‘without supervised labels of ID’ in Table 3, we train MoCo-v2 on CIFAR-10. The model is trained over 800 epochs with the SGD optimizer and its momentum 0.9. The batch size is 512. The learning rate is decayed by the cosine scheduler from 0.06 to 0. The model backbone is ResNet-18 combined with an MLP projection head. For the other configurations, we follow those given in the link¹. After training the MoCo-v2 model, the NAN score is computed over multiple (9 overall) translated images of the test sample including the original image, and the scores are aggregated by average [43]. This aggregation technique is used exclusively for the model trained by MoCo-v2.

Setup: OOD CIFAR-10 The model training configuration for OCC is similar to that of label-free OOD detection on CIFAR-10 except that the train dataset is augmented randomly with 90-degree rotations. During the inference, the rotation is not used.

C.2. Score Fusion

A distance-based score $S_{dist}(\mathbf{x}) = d(X_{ind}, \mathbf{x})$ (e.g. KNN, SSD, or Mahalanobis) can be combined with NAN in a simple manner by

$$S_{dist+NAN}(\mathbf{x}) = d(X_{ind}, \mathbf{x}) / \|\mathbf{a}^{(L)}\|_{NAN}. \quad (49)$$

¹https://colab.research.google.com/github/facebookresearch/moco/blob/colab-notebook/colab/moco_cifar10_demo.ipynb

ID	Architecture	Last hidden layer $\mathbf{a}^{(L)}$	AUROC \uparrow		FPR95 \downarrow	
			l_1 -norm / NAN		l_1 -norm / NAN	
CIFAR-10	ResNet-18	average pool	93.27 / 93.56 (+ 0.29)		40.42 / 38.86 (- 1.56)	
	ResNet-18 + projection head	hidden layer in projection head	92.43 / 94.94 (+ 2.51)		43.02 / 30.08 (- 12.94)	
ImageNet-1k	ResNet-50	average pool	87.09 / 86.33 (- 0.76)		44.67 / 46.56 (+ 1.89)	
	ResNet-50 + projection head	hidden layer in projection head	57.99 / 92.32 (+ 34.33)		95.22 / 31.59 (- 63.63)	

Table 6: **Ablation of NAN with respect to the projection head.** The sparsity term in NAN is particularly effective when applied to the network architecture that contains the MLP projection head. Note that the l_1 -norm here refers to the NAN score without the sparsity term. The reported performance here is obtained by averaging over all test OOD datasets.

	ReLU		Leaky ReLU		GeLU	
	AUROC \uparrow	FPR95 \downarrow	AUROC \uparrow	FPR95 \downarrow	AUROC \uparrow	FPR95 \downarrow
NAN w/o sparsity term (l_1 -norm)	92.43	43.02	92.40	44.65	92.68	43.84
NAN	94.94	30.08	94.92	30.56	94.05	35.02

Table 7: **Ablation of NAN with respect to the activation functions** used in the last hidden layer. The ID data is CIFAR-10. The results indicate two aspects: (1) The performance of NAN is fairly robust with different choices of the activation function. (2) The sparsity term in NAN is always effective. The reported performance here is obtained by averaging over all test OOD datasets.

ID	Formula	ImageNet-1k			CIFAR-10		
		d	AUROC \uparrow	FPR95 \downarrow	d	AUROC \uparrow	FPR95 \downarrow
embedding magnitude	$\ g(\mathbf{x})\ _2$	128	84.09	72.85	128	93.00	43.40
NAN w/o sparsity term	$\ \mathbf{a}^{(L)}\ _1$	2048	57.99	95.22	512	92.40	43.00
NAN	$\ \mathbf{a}^{(L)}\ _{\text{NAN}}$	2048	92.32	31.59	512	94.90	30.10

Table 8: Comparison of NAN with the embedding magnitude. The embedding magnitude has been widely used in previous works. Here d indicates the dimension of the corresponding layer. The dimension of the embedding layer is often chosen small for effective training of the model. Due to its small layer dimension, the embedding magnitude may not fully capture the activation patterns, and hence can be sub-optimal. The reported performance here is obtained by averaging over all test OOD datasets.

D. Further Analysis on NAN

Setup We follow the same setup given in Sec. 6. When CIFAR-10 is the ID data, the test OOD datasets are LSUN-fix, ImageNet-fix, CIFAR-100, SVHN, and Places. When ImageNet-1k is the ID data, the test OOD datasets are iNaturalist, SUN, Places, and Texture.

D.1. Analysis on Projection Head

We analyze NAN with respect to **the projection head**. Table 6 indicates that NAN is more effective when it is applied to the hidden layer of the projection head rather than the average pooling layer.

NAN (*i.e.* particularly its sparsity term) becomes effective when the network learns to increase the number of deactivated units of ID samples (or have a relatively larger number of deactivated units for ID samples than OOD instances). Due to the entanglement of the feature map units in the average pooling layer, the network may not effectively increase the number of deactivated units in the average pooling layer. Hence, NAN can be sub-optimal for the average pooling layer.

D.2. Analysis on Activation Function

We evaluate NAN with **different activation functions**. We follow the same experimental protocol given in Sec. 6.3. We apply different activation functions in the hidden layer of the projection head. The results given in Table 7 shows that NAN is robust with respect to the choice of the activation function.

D.3. Comparison with Embedding Magnitude

For the sake of extensiveness, we compare NAN with the **embedding magnitude**. The embedding magnitude has been widely used in prior works for OOD detection-related tasks. The dimension of the embedding layer is often chosen to be a small number to avoid the curse of dimensionality during training. This may have a trade-off to OOD detection as the

OOD	iNaturalist		SUN		Places		Texture		Average		ID ACC
	AUROC \uparrow	FPR95 \downarrow									
MSP	89.63	50.57	80.64	75.54	79.78	76.24	82.98	65.14	83.26	66.87	81.07
Energy	83.76	49.68	56.50	75.22	54.77	78.38	72.44	65.09	66.87	67.09	81.07
Mahalanobis	91.96	43.76	75.62	86.01	61.50	89.74	84.60	67.93	78.42	71.86	81.07
KNN	91.43	50.04	83.45	75.76	79.46	78.41	89.25	50.78	85.90	63.75	81.07
embedding magnitude	81.26	66.16	78.64	67.44	75.81	69.37	82.93	57.11	79.66	65.02	81.07
NAN w/o sparsity term (<i>i.e.</i> l_1 -norm)	54.93	83.98	67.05	80.47	65.25	81.01	67.87	72.54	63.78	79.50	81.07
NAN	92.46	45.82	82.11	67.62	80.46	69.66	87.24	57.77	85.57	60.22	81.07

Table 9: Results on ImageNet-1k (ID) with ViT-B/16.

test OOD datasets	Formula	LSUN-fix		ImageNet-fix		CIFAR-100		SVHN		Places		Average	
		AUROC \uparrow	FPR95 \downarrow										
hidden classifier confidence	$\max_k \overline{\psi}_k^{(L)}(\mathbf{x})$	95.06	33.35	94.54	35.92	92.17	45.10	94.66	39.91	94.66	30.15	94.22	36.89

Table 10: Results on CIFAR-10 (ID) with ResNet-18. The hidden classifier confidence is evaluated as a score function for OOD detection. The results shows that the hidden classifier confidence is capable of OOD detection.

embedding of a small dimension may not capture diverse activation patterns of embedding layer units and therefore its norm may not effectively differentiate OOD from ID. This hypothesis seems consistent to the results given in Table 8.

D.4. Evaluation of NAN on ViT

We evaluate NAN on the **vision transformer ViT**. We utilize ViT-B/16 pretrained on ImageNet-1k, which can be downloaded from PyTorch². Analogous to the observations in Sec. D.1, direct usage of NAN on the pretrained ViT can be sub-optimal because the class token output of ViT is the LayerNorm layer, which can cancel out the norm information therein. Therefore, we add an MLP projection head on top of the pretrained ViT, and fine-tune the projection head while freezing the pretrained ViT backbone. The MLP projection head consists of a single hidden layer whose dimension is 786 and its activation function is ReLU. The embedding of the projection head is normalized and divided by the temperature 0.2, and trained by the cross entropy with 10 epochs under SGD, using the learning rate 0.03 that decays to 0 by the cosine scheduler.

For comparison, the KNN and Mahalanobis scores are applied on the original class token output of the pretrained ViT, and hence are independent of the projection head fine-tuning. Other OOD detection scores (MSP, Energy, and embedding magnitude) are applied to the fine-tuned classifier of the projection head. NAN utilizes the hidden layer in the projection head as this layer is the last hidden layer that involves the activation function computation.

Table 9 shows that NAN is effective for the ViT network as well. In addition, NAN is comparable to the state-of-the-art OOD detection scores.

Note on the ViT performance of KNN Note that the performance of KNN in Table 9 is lower than that of KNN reported in [42]. This is because the KNN we implemented is applied on ViT pretrained on ImageNet-1k, while the KNN reported in [42] is applied on ViT pretrained on ImageNet-21k.

D.5. Evaluation of Hidden Classifier for OOD Detection

We evaluate the **hidden classifier for OOD detection**. NAN’s numerator is the l_1 -norm of the activation vector, which we proved is a confidence value of the hidden classifier. We test this numerator component by testing the OOD detection capability of this hidden classifier confidence. Table 10 shows the hidden classifier confidence is capable of OOD detection.

D.6. Evaluation of NAN on CIDER

CIDER [29] is a training framework that is particularly effective for the KNN score. We evaluate NAN’s compatibility to the KNN score from the model trained by CIDER. The results shown in Table 11 indicates that NAN can effectively enhance the KNN score of CIDER.

	SVHN		Places365		iSUN		Texture		LSUN		Average	
	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC
NAN	73.82	90.46	26.33	94.65	25.47	96.46	25.35	95.21	1.17	99.45	30.43	95.25
KNN	4.44	99.36	37.88	92.97	22.94	96.16	17.27	97.15	9.85	98.21	18.48	96.77
NAN+KNN	5.70	98.62	21.79	95.32	14.01	97.64	16.21	96.61	0.95	99.68	11.73	97.57

Table 11: The results of the OOD detection scores (KNN, NAN, NAN+KNN) on the model trained by CIDER on CIFAR-10 (ID).

	iNaturalist		SUN		Places		Texture		ImageNet-O		OpenImage-O		Species		Average	
	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC
l_1 -norm	97.52	52.06	95.58	59.40	95.65	61.30	92.11	59.21	88.20	67.97	92.43	63.10	95.83	59.42	93.90	60.35
$1/l_0$ -norm	15.66	96.58	33.38	91.83	39.10	90.37	44.36	87.41	88.60	56.76	41.29	88.58	64.04	79.55	46.63	84.44
Residual	28.74	95.09	46.88	89.76	58.91	85.77	11.28	96.45	63.50	84.24	34.96	93.34	74.43	73.72	45.53	88.34
NAN	15.86	96.94	29.81	92.77	37.21	91.46	43.46	88.09	87.95	69.74	38.12	92.44	64.56	80.09	45.28	87.36
<i>with ReAct:</i>																
l_1 -norm	98.07	37.19	96.37	46.97	96.90	45.47	85.44	61.21	84.95	74.80	93.54	54.48	98.81	41.25	93.44	51.62
$1/l_0$ -norm	21.19	95.60	36.56	90.81	41.28	89.63	52.16	82.23	90.35	53.37	49.25	85.85	61.45	81.89	50.32	82.77
Residual	28.59	95.06	39.40	91.95	51.02	88.18	12.11	96.87	68.30	83.01	36.67	92.62	72.27	75.03	44.05	88.96
NAN	13.86	97.37	24.90	94.69	33.31	92.52	34.02	91.44	84.10	71.72	37.27	92.02	63.68	81.10	41.59	88.69

Table 12: The comparison of NAN with various forms of vector norms on ImageNet-1k (ID).

D.7. Comparison of NAN to various forms of vector norms

To further highlight the effectiveness of NAN, we compare NAN with various forms of vectors norms; namely, l_1 -norm, the reciprocal of l_0 -norm, and the residual of ViM which is the l_2 -norm of the orthogonal projection. The experiment protocol follows [51], and the OOD datasets can be downloaded from its GitHub repository.

The results in Table 12 indicate that NAN is significantly better than the l_1 -norm and the reciprocal of l_0 -norm. We note that l_1 -norm does not capture deactivation, while the reciprocal of l_0 -norm captures only deactivation. Hence, the superiority of NAN over these vector norms indicate that capturing both activation and deactivation is crucial.

Compared to the residual of ViM, on the other hand, NAN is notably superior with respect to the FPR95 metric when ReAct is applied on the model, while NAN is comparable to the residual when without ReAct. We note, however, that the residual of ViM requires eigen decomposition of the bankset features, while the computation of NAN is done by a single forward pass of the network.

E. Limitation of NAN

Based on our theoretical observations, NAN is intrinsically a classifier output and hence may inherit the weaknesses of classifier-based OOD detectors that have been recently found in [8, 9]. In addition, as observed in Sec. D.1, the optimal usage of NAN requires networks that involve the MLP projection head.

²https://pytorch.org/vision/main/models/generated/torchvision.models.vit_b_16.html

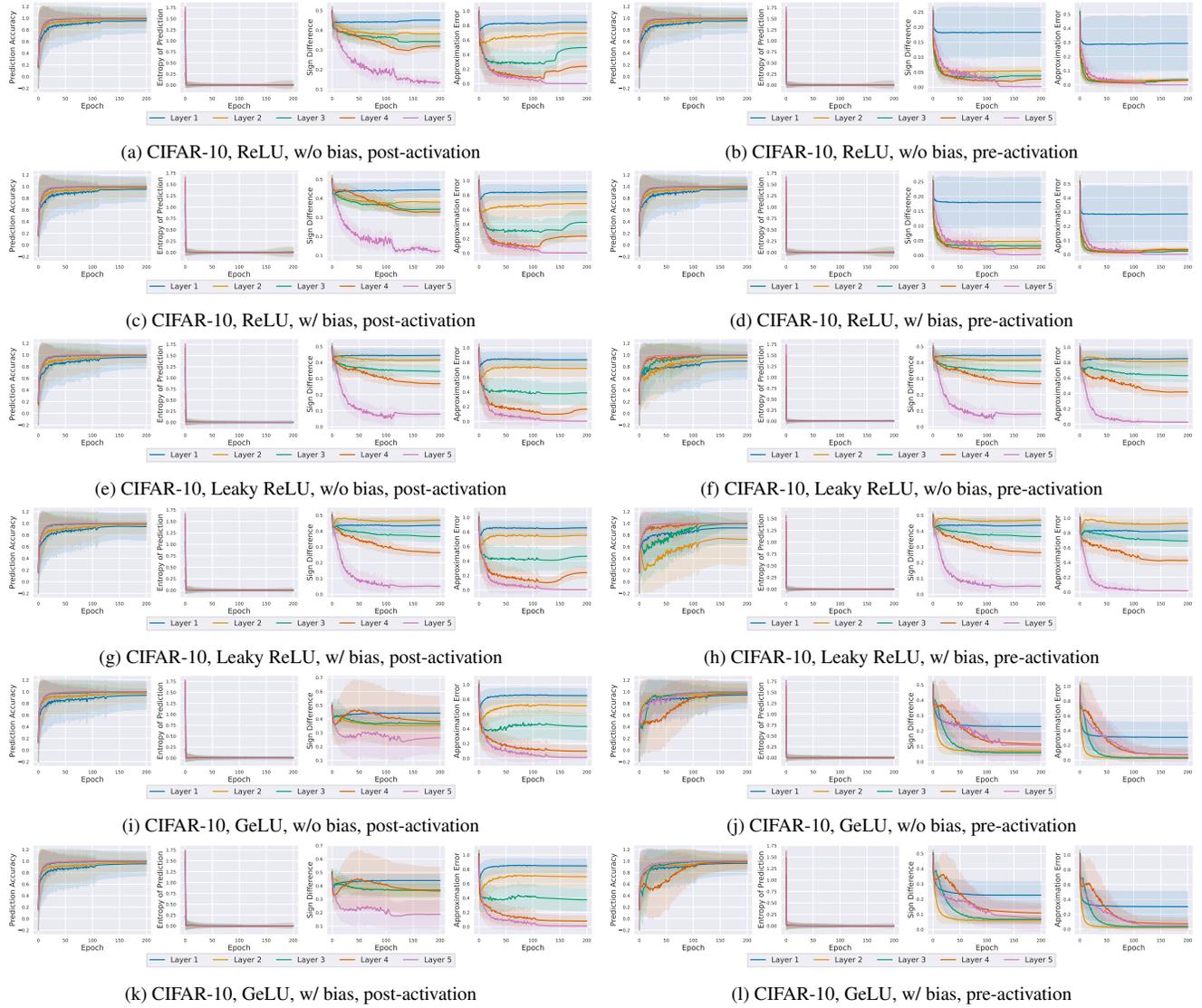


Figure 12: Results of hidden classifiers with different activation functions (ReLU, Leaky ReLU, and GeLU) on CIFAR-10.

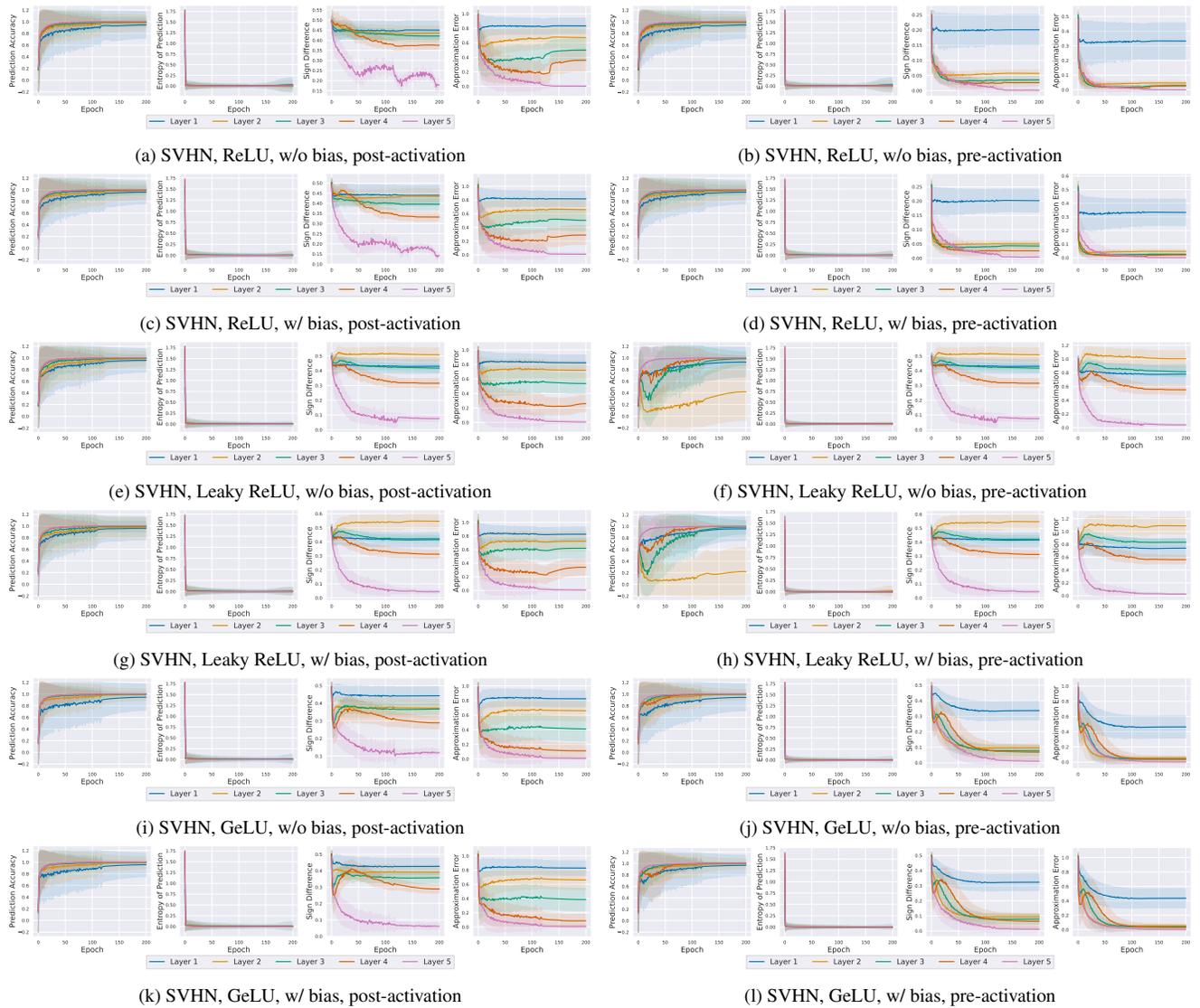


Figure 13: Results of hidden classifiers with different activation functions (ReLU, Leaky ReLU, and GeLU) on SVHN.

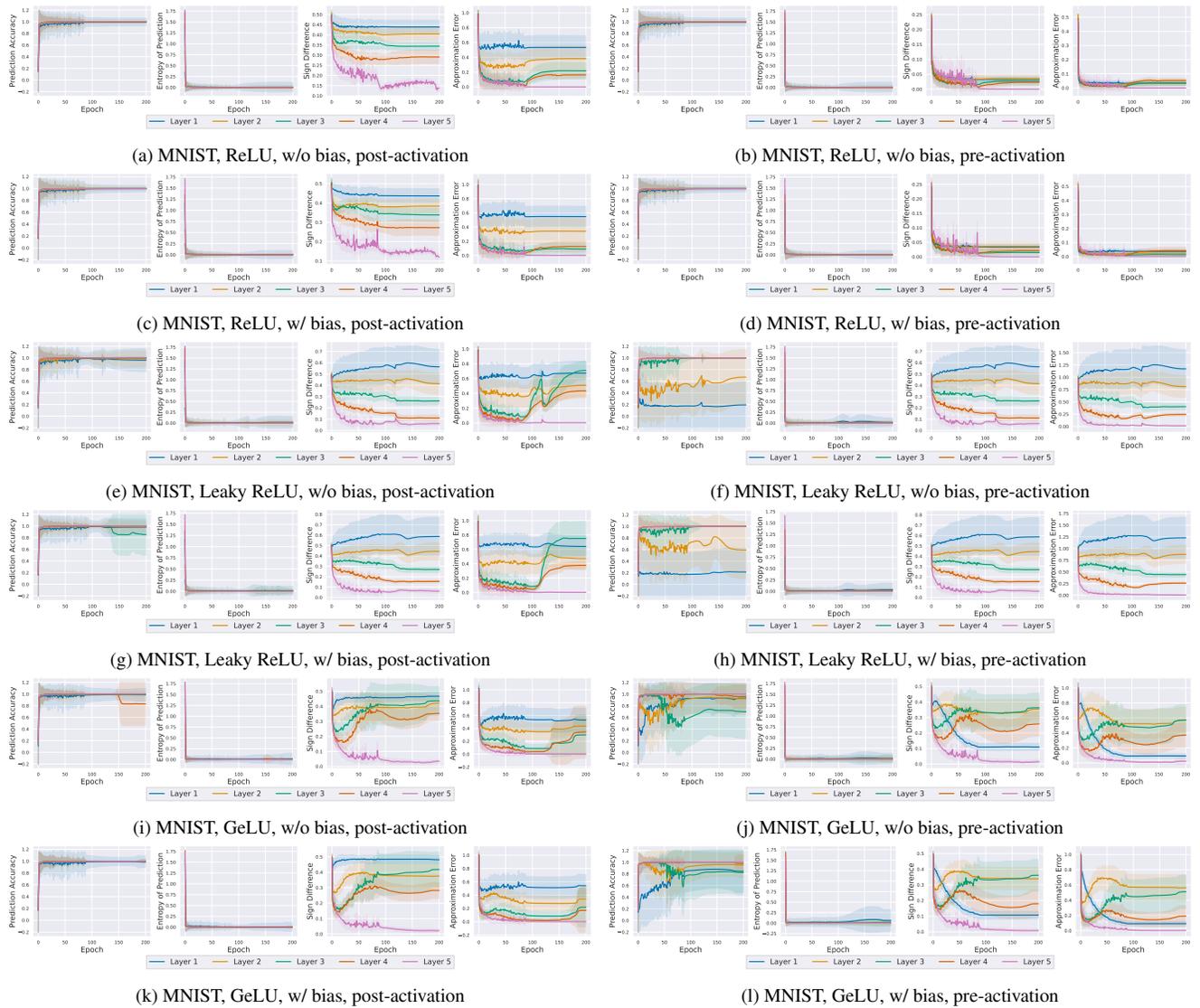


Figure 14: Results of hidden classifiers with different activation functions (ReLU, Leaky ReLU, and GeLU) on MNIST.