

# Full-Spectrum Out-of-Distribution Detection

Jingkang Yang · Kaiyang Zhou · Ziwei Liu

Received: date / Accepted: date

**Abstract** Existing out-of-distribution (OOD) detection literature clearly defines semantic shift as a sign of OOD but does not have a consensus over covariate shift. Samples experiencing covariate shift but not semantic shift are either excluded from the test set or treated as OOD, which contradicts the primary goal in machine learning—being able to generalize beyond the training distribution. In this paper, we take into account both shift types and introduce full-spectrum OOD (FS-OOD) detection, a more realistic problem setting that considers both detecting semantic shift and being tolerant to covariate shift; and designs three benchmarks. These new benchmarks have a more fine-grained categorization of distributions (*i.e.*, training ID, covariate-shifted ID, near-OOD, and far-OOD) for the purpose of more comprehensively evaluating the pros and cons of algorithms. To address the FS-OOD detection problem, we propose SEM, a simple feature-based semantics score function. SEM is mainly composed of two probability measures: one is based on high-level features containing both semantic and non-semantic information, while the other is based on low-level feature statistics only capturing non-semantic image styles. With a simple combination, the non-semantic part is cancelled out, which leaves only semantic information in SEM that can better handle FS-OOD detection. Extensive experiments on the three new benchmarks show

that SEM significantly outperforms current state-of-the-art methods. Our code and benchmarks are released in <https://github.com/Jingkang50/OpenOOD>.

## 1 Introduction

State-of-the-art deep neural networks are notorious for their overconfident predictions on out-of-distribution (OOD) data [1], defined as those not belonging to in-distribution (ID) classes. Such a behavior makes real-world deployments of neural network models untrustworthy and could endanger users involved in the systems. To solve the problem, various OOD detection methods have been proposed in the past few years [2, 3, 4, 5, 6, 7, 8]. The main idea for an OOD detection algorithm is to assign to each test image a score that can represent the likelihood of whether the image comes from in- or out-of-distribution. Images whose scores fail to pass a threshold are rejected, and the decision-making process should be transferred to humans for better handling.

A critical problem in existing research of OOD detection is that only semantic shift is considered in the detection benchmarks while covariate shift—a type of distribution shift that is mainly concerned with changes in appearances like image contrast, lighting or viewpoint—is either excluded from the evaluation stage or simply treated as a sign of OOD [1], which contradicts with the primary goal in machine learning, *i.e.*, to generalize beyond the training distribution [9].

In this paper, we introduce a more challenging yet realistic problem setting called *full-spectrum out-of-distribution detection*, or *FS-OOD detection*. The new setting takes into account both the detection of semantic shift and the ability to recognize covariate-shifted data as ID. To this end, we design three benchmarks,

---

Jingkang Yang  
S-Lab, Nanyang Technological University, Singapore  
E-mail: jingkang001@ntu.edu.sg

Kaiyang Zhou  
S-Lab, Nanyang Technological University, Singapore  
E-mail: kaiyang.zhou@ntu.edu.sg

Ziwei Liu  
S-Lab, Nanyang Technological University, Singapore  
E-mail: ziwei.liu@ntu.edu.sg

namely DIGITS, OBJECTS and COVID, each targeting a specific visual recognition task and together constituting a comprehensive testbed. We also provide a more fine-grained categorization of distributions for the purpose of thoroughly evaluating an algorithm. Specifically, we divide distributions into four groups: training ID, covariate-shifted ID, near-OOD, and far-OOD (the latter two are inspired by a recent study [10]). Figure 1-a shows example images from the DIGITS benchmark: the covariate-shifted images contain the same semantics as the training images, *i.e.*, digits from 0 to 9, and should be classified as ID, whereas the two OOD groups clearly differ in semantics but represent two different levels of covariate shift.

Ideally, an OOD detection system is expected to produce high scores for samples from the training ID and covariate-shifted ID groups, while assign low scores to samples from the two OOD groups. However, when applying a state-of-the-art OOD detection method, *e.g.* the energy-based EBO [4], to the proposed benchmarks like DIGITS (see Figure 1-b), we observe that the resulting scores completely fail to distinguish between ID and OOD. As shown in Figure 1-b, all data are classified as ID including both near-OOD and far-OOD samples.

To address the more challenging but realistic FS-OOD detection problem, we propose SEM, *a simple feature-based semantics score function*. Unlike existing score functions that are based on either marginal distribution [4] or predictive confidence [2], SEM leverages features from both top and shallow layers to deduce a single score that is only relevant to semantics, hence more suitable for identifying semantic shift while ensuring robustness under covariate shift. Specifically, SEM is mainly composed of two probability measures: one is based on high-level features containing both semantic and non-semantic information, while the other is based on low-level feature statistics only capturing non-semantic image styles. With a simple combination, the non-semantic part is cancelled out, which leaves only semantic information in SEM. Figure 1-c illustrates that SEM’s scores are much clearer to distinguish between ID and OOD.

We summarize the **contributions** of this paper as follows. **1)** For the first time, we introduce the full-spectrum OOD detection problem, which represents a more realistic scenario considering both semantic and covariate shift in the evaluation pipeline. **2)** Three benchmark datasets are designed for research of FS-OOD detection. They cover a diverse set of recognition tasks and have a detailed categorization over distributions. **3)** A simple yet effective OOD detection score function called SEM is proposed. Through extensive experiments on the three new benchmarks, we

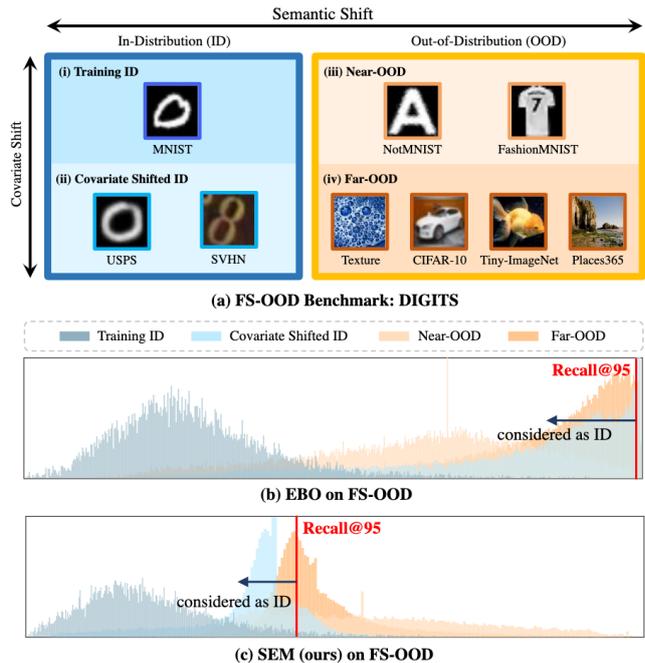


Fig. 1: **Comparison of OOD detection scores obtained by different approaches on the newly introduced full-spectrum OOD detection benchmark: (a) DIGITS Benchmark.** Ideally, the scores should be clear enough to separate out OOD data while include covariate-shifted data as in-distribution—which has been ignored by most existing research. (b) The state-of-the-art energy-based approach, EBO [4], apparently fails in this scenario. (c) Our approach, based on a semantics-oriented score function, can improve the detection performance significantly.

demonstrate that SEM significantly outperforms current state-of-the-art methods in FS-OOD detection. The source code and new datasets are open-sourced in <https://github.com/Jingkang50/OpenOOD>.

## 2 Related Work

The key idea in out-of-distribution (OOD) detection is to design a metric, known as score function, to assess whether a test sample comes from in- or out-of-distribution. The most commonly used metric is based on the conditional probability  $p(y|x)$ . An early OOD detection method is maximum softmax probability (MSP) [2], which is motivated by the observation that deep neural networks tend to give lower confidence to mis-classified or OOD data. A follow-up work ODIN [3] applies a temperature scaling parameter to soften the probability distribution, and further improves the performance by injecting adversarial perturbations to the

input. Model ensembling has also been found effective in enhancing robustness in OOD detection [11, 12].

Another direction is to design the metric in a way that it reflects the marginal probability  $p(\mathbf{x})$ . Liu *et al.* [4] connect their OOD score to the marginal distribution using an energy-based formulation, which essentially sums up the prediction logits over all classes. Lee *et al.* [5] assume the source data follow a normal distribution and learn a Mahalanobis distance to compute the discrepancy between test images and the estimated distribution parameters. Generative modeling has also been investigated to estimate a likelihood ratio for scoring test images [11, 6, 13].

Some methods exploit external OOD datasets. For example, Hendrycks *et al.* [14] extend MSP by training the model to produce uniform distributions on external OOD data. Later works introduce re-sampling strategy [15] and cluster-based methodology [16] to better leverage the background data. However, this work do not use external OOD datasets for model design.

Different from all existing methods, our approach aims to address a more challenging scenario, *i.e.*, FS-OOD detection, which has not been investigated in the literature but is critical to real-world applications. The experiments show that current state-of-the-art methods mostly fail in the new setting while our approach gains significant improvements.

### 3 Methodology

#### 3.1 Feature-Based Semantics Score Function

Key to detect out-of-distribution (OOD) data lies in the design of a score function, which is used as a quantitative measure to distinguish between in- and out-of-distribution data. Our idea is to design the function in such a way that the degree of semantic shift is effectively captured, *i.e.*, the designed score to be only sensitive to semantic shift while being robust to covariate shift. For data belonging to the in-distribution classes, the score is high, and vice versa.

**Formulation** Our score function, called SEM, has the following design:

$$\text{SEM}(\mathbf{x}) = \log p(\mathbf{x}_s), \quad (1)$$

where  $\mathbf{x}$  denotes image features learned by a neural network; and  $\mathbf{x}_s$  denotes features that only capture the semantics. The probability  $p(\mathbf{x}_s)$  can be computed by a probabilistic model, such as a Gaussian mixture model.

The straightforward way to model  $\mathbf{x}_s$  is to learn a neural network for image recognition and hope that the output features  $\mathbf{x}$  only contain semantic information,

*i.e.*,  $\mathbf{x}_s = \mathbf{x}$ . If so, the score can be simply computed by  $\text{SEM}(\mathbf{x}) = \log p(\mathbf{x})$ . However, numerous studies have suggested that the output features  $\mathbf{x}$  often contain both semantic and non-semantic information while decoupling them is still an open research problem [9, 18, 19]. Let  $\mathbf{x}_n$  denote non-semantic features, we assume that semantic features  $\mathbf{x}_s$  and non-semantic features  $\mathbf{x}_n$  are generated independently, namely

$$p(\mathbf{x}) = p(\mathbf{x}_s)p(\mathbf{x}_n). \quad (2)$$

We propose a simple method to model the score function so that it becomes only relevant to the semantics of an image. This is achieved by leveraging *low-level feature statistics*, *i.e.*, means and standard deviations, learned in a CNN, which have been shown effective in capturing image styles that are essentially irrelevant to semantics [20]. Specifically, the score function in Eq. 1 is rewritten as

$$\text{SEM}(\mathbf{x}) = \log p(\mathbf{x}_s) = \log \frac{p(\mathbf{x}_s)p(\mathbf{x}_n)}{p(\mathbf{x}_n)} = \log \frac{p(\mathbf{x})}{p(\mathbf{x}_n)}, \quad (3)$$

where  $p(\mathbf{x})$  is computed using the output features while  $p(\mathbf{x}_n)$  is based on low-level feature statistics.

Below we first discuss how to compute feature statistics and then detail the approach of how to model the distributions for  $\mathbf{x}$  and  $\mathbf{x}_n$ .

**Feature Statistics Computation** Instance-level feature statistics have been widely used in the style transfer community for manipulating image style [21]. Given a set of CNN feature maps  $\mathbf{z} \in \mathbb{R}^{C \times H \times W}$  with  $C$ ,  $H$  and  $W$  denoting the number of channels, height and width, their feature statistics, *i.e.*, means  $\boldsymbol{\mu} \in \mathbb{R}^C$  and standard deviations  $\boldsymbol{\sigma} \in \mathbb{R}^C$ , are computed across the spatial dimension within each channel  $c = \{1, 2, \dots, C\}$ ,

$$\mu_c = \frac{1}{HW} \sum_{h=1}^H \sum_{w=1}^W z_{c,h,w}, \quad (4)$$

$$\sigma_c = \left( \frac{1}{HW} \sum_{h=1}^H \sum_{w=1}^W (z_{c,h,w} - \mu_c)^2 \right)^{\frac{1}{2}}. \quad (5)$$

As shown in Zhou *et al.* [20], the feature statistics in shallow CNN layers are strongly correlated with domain information (*i.e.*, image style) while those in higher layers pick up more semantics. Therefore, we choose to extract feature statistics in the first CNN layer and represent  $\mathbf{x}_n$  by concatenating the means and standard deviations, *i.e.*,  $\mathbf{x}_n = [\boldsymbol{\mu}, \boldsymbol{\sigma}]^T$ .

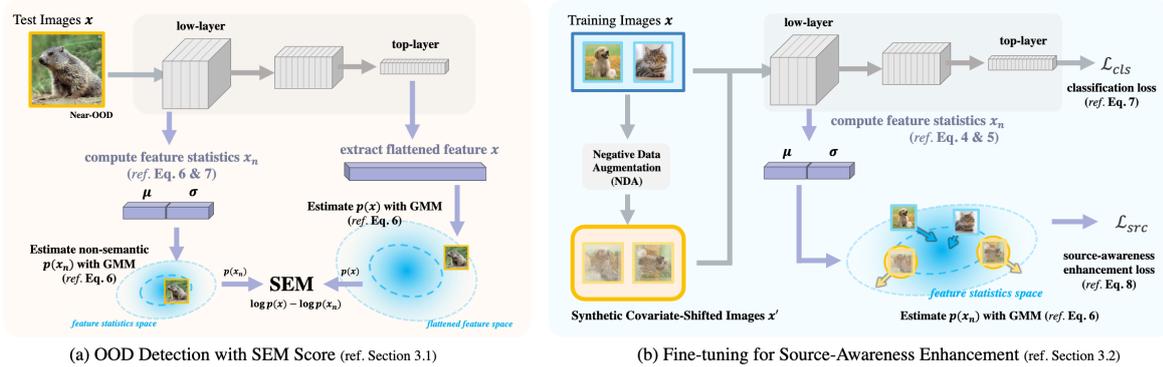


Fig. 2: **Overview of our Methodology.** (a) The computation of SEM score function for OOD detection. SEM combines the estimation of  $p(\mathbf{x})$  (using top-layer features to capture both semantic and non-semantic information) and  $p(\mathbf{x}_n)$  (using low-level feature statistics to only capture non-semantic information) with Eq. 3 for better concentration on semantics. (b) The fine-tuning scheme to enhance source-awareness for better estimating  $p(\mathbf{x}_n)$ . The main idea is to pull together the instance-level CNN feature statistics of in-distribution data to make them more compact, while pushing away those of synthetic OOD data, which are obtained by negative data augmentation such as Mixup [17].

**Distribution Modeling** For simplicity, we model  $p(\mathbf{x})$  and  $p(\mathbf{x}_n)$  in Eq. 3 using the same approach, which consists of two steps: dimension reduction and distribution modeling. Below we only discuss  $p(\mathbf{x})$  for clarity.

Motivated by the manifold assumption in Bengio *et al.* [22] that suggests data typically lie in a manifold of much lower dimension than the input space, we transform features  $\mathbf{x}$  to a new low-dimensional space, with a hope that the structure makes it easier to distinguish between in- and out-of-distribution. To this end, we propose a variant of the principal component analysis (PCA) approach. Specifically, rather than maximizing the variance for the entire population, we maximize the sum of variances computed within each class with respect to the transformation matrix. In doing so, we can identify a space that is less correlated with classes.

Given a training dataset, we build a Gaussian mixture model (GMM) to capture  $p(\mathbf{x})$ . Formally,  $p(\mathbf{x})$  is defined as

$$p(\mathbf{x}) = \sum_{m=1}^M \lambda_m \mathcal{N}(\boldsymbol{\alpha}_m, \boldsymbol{\beta}_m), \quad (6)$$

where  $M$  denotes the number of mixture components,  $\lambda_m$  the mixture weight s.t.  $\sum_{m=1}^M \lambda_m = 1$ , and  $\boldsymbol{\alpha}_m$  and  $\boldsymbol{\beta}_m$  the means and variances of a normal distribution. A GMM model can be efficiently trained by the expectation-maximization (EM) algorithm.

### 3.2 Source-Awareness Enhancement

While feature statistics exhibit a higher correlation with source distributions [20], the boundary between in- and

out-of-distribution in complicated real-world data is not guaranteed to be clear enough for differentiation. Inspired by Liu *et al.* [4] who fine-tune a pretrained model to increase the energy values assigned to OOD data and lower down those for ID data, we propose a fine-tuning scheme to enhance source-awareness in feature statistics. An overview of the fine-tuning scheme is illustrated in Figure 2-b.

**Negative Data Augmentation** The motivation behind our fine-tuning scheme is to obtain a better estimate of non-semantic score, in hope that it will help SEM better capture the semantics with the combination in Eq. 3. This can be achieved by explicitly training feature statistics of ID data to become more compact, while pushing OOD data’s feature statistics away from the ID support areas. A straightforward way is to collect auxiliary OOD data like Liu *et al.* [4] for building a contrastive objective. In this work, we propose a more efficient way by using negative data augmentation [23] to synthesize OOD samples. The key idea is to choose data augmentation methods to easily generate samples with covariate shift. One example augmentation is Mixup [17].

**Learning Objectives** Given a source dataset  $\mathcal{S} = \{(\mathbf{x}, y)\}$ ,<sup>1</sup> we employ negative data augmentation methods  $\text{aug}(\cdot)$  to synthesize an OOD dataset  $\mathcal{S}_{aug} = \{(\mathbf{x}', y)\}$  where  $\mathbf{x}' = \text{aug}(\mathbf{x})$ . For fine-tuning, we combine a classification loss  $\mathcal{L}_{cls}$  with a source-awareness enhancement loss  $\mathcal{L}_{src}$ . These two losses are formally

<sup>1</sup> With a slight abuse of notation, we use  $\mathbf{x}$  here to denote an image.

defined as

$$\mathcal{L}_{cls} = - \sum_{(\mathbf{x}, y) \sim \mathcal{S}} \log p(y|\mathbf{x}), \quad (7)$$

and

$$\mathcal{L}_{src} = \sum_{\mathbf{x}' \sim \mathcal{S}_{aug}} p(\mathbf{x}'_n) - \sum_{\mathbf{x} \sim \mathcal{S}} p(\mathbf{x}_n), \quad (8)$$

where the marginal probability  $p(\mathbf{x})$  is computed based on a GMM model described previously. Note that the GMM model is updated every epoch to adapt to the changing features.

After fine-tuning, we learn a new GMM model using the original source dataset. This model is then used to estimate the marginal probability  $p(\mathbf{x})$  at test time.

#### 4 FS-OOD Benchmarks

To evaluate full-spectrum out-of-distribution (FS-OOD) detection algorithms, we design three benchmarks: DIGITS, OBJECTS, and COVID. Examples for DIGITS are shown in Figure 1 and the other two are shown in Figure 3.

**Benchmark-1: DIGITS** We construct the DIGITS benchmark based on the popular digit datasets: MNIST [24], which contains 60,000 images for training. During testing, the model will be exposed to 10,000 MNIST test images, with 26,032 covariate-shifted ID images from SVHN [25] and another 9,298 from USPS [26]. The near-OOD datasets are notMNIST [27] and FashionMNIST [28], which share a similar background style with MNIST. The far-OOD datasets consist of a textural dataset (Texture [29]), two object datasets (CIFAR-10 [30] & Tiny-ImageNet [31]), and one scene dataset (Places365 [32]). The CIFAR-10 and Tiny-ImageNet test sets have 10,000 images for each. The Places365 test set contains 36,500 scene images.

**Benchmark-2: OBJECTS** The OBJECTS benchmark is built on top of CIFAR-10 [30], which contains 50,000 images for training. During testing, the model will be exposed to 10,000 CIFAR-10 test images, and another 10,000 images selected from ImageNet-22K [31] with the same categories as CIFAR-10 (so it is called ImageNet-10). For ImageNet-10, we choose five ImageNet-22K classes corresponding to one CIFAR-10 class, with each class selecting 1,000 training images and 200 testing images. Details of the selected classes are shown in Table 1. In addition to ImageNet, CIFAR-10-C is used as a covariate-shifted ID dataset, which is essentially a corrupted version of CIFAR-10. For near-OOD, we choose CIFAR-100 and Tiny-ImageNet. For far-OOD, we choose MNIST, FashionMNIST, Texture and CIFAR-100-C.

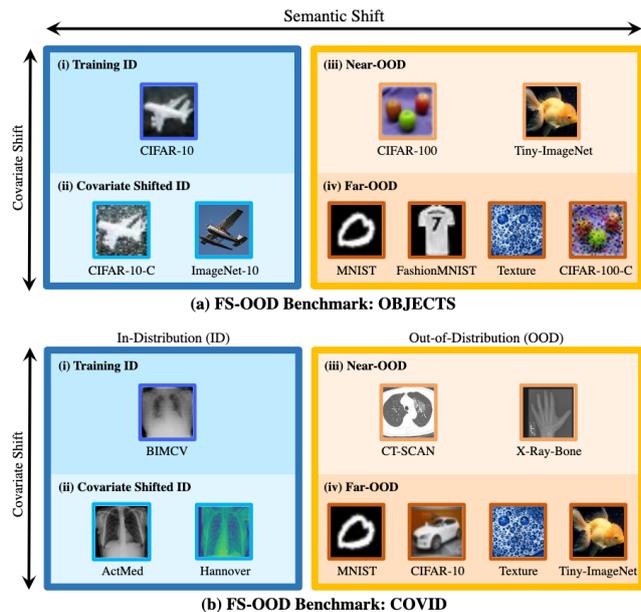


Fig. 3: Examples for the two FS-OOD detection benchmarks: COVID and OBJECTS. Each benchmark consists of a training ID dataset, two covariate-shifted ID datasets, two near-OOD datasets, and four far-OOD datasets.

**Benchmark-3: COVID** We construct a real-world benchmark to show the practical value of FS-OOD. We simulate the scenario where an AI-assisted diagnostic system is trained to identify COVID-19 infection from chest x-ray images. The training data come from a single source (*e.g.*, a hospital) while the covariate-shifted ID test data are from other hospitals or machines, to which the system needs to be robust and produce reliable predictions. Specifically, we refer to the COVID-19 chest X-ray dataset review [33], and use the large-scale image collection from Valencian Region Medical Image-Bank [34] (referred to as BIMCV) as training ID images (randomly sampled 2443 positive cases and 2501 negative cases with necessary cleaning). Images from two other sources, *i.e.*, ACTUALMED [35] (referred to as ActMed with 132 positive images), and Hannover [36] (from Hannover Medical School with 243 positive images), are considered as the covariate-shifted ID group. OOD images are from completely different classes. Near-OOD images are obtained from other medical datasets, *i.e.*, the RSNA Bone Age dataset with 200 bone X-ray images [37] and 544 COVID CT images [38]. Far-OOD samples are defined as those with drastic visual and concept differences than the ID images. We use MNIST, CIFAR-10, Texture and Tiny-ImageNet.

**Evaluation Metrics** In the FS-OOD setting, different datasets belonging to one OOD type (*i.e.*, near-OOD or far-OOD) are grouped together. We also re-

Table 1: **Selected ImageNet-22K classes for OBJECTS benchmark.** We manually find 5 ImageNet-22K classes that belong to each CIFAR-10 classes, and pick the first 1,000 images from every selected class for OBJECTS benchmark. A string such as ‘n03365231’ is the synset id for downloading the corresponding class from ImageNet API.

Airplane	Automobile	Bird	Cat
n03365231 floatplane	n04516354 used car	n01503061 bird	n02121808 domestic cat
n02691156 airplane	n04285008 sports car	n01812337 dove	n02123159 tiger cat
n04552348 warplane	n02958343 car	n01562265 robin	n02122878 tabby
n02686568 aircraft	n03594945 jeep	n01539573 sparrow	n02123394 Persian cat
n02690373 airliner	n02930766 cab	n01558594 blackbird	n02123597 Siamese cat

Deer	Dog	Frog
n02430045 deer	n02116738 African hunting dog	n01639765 frog
n02431122 red deer	n02087122 hunting dog	n01641577 bullfrog
n02432511 mule deer	n02105855 Shetland sheepdog	n01644373 tree frog
n02433318 fallow deer	n02109961 Eskimo dog	n01640846 true frog
n02431976 Japanese deer	n02099601 golden retriever	n01642539 grass frog

Horse	Ship	Truck
n02387254 farm horse	n02965300 cargo ship	n04490091 truck
n02381460 wild horse	n04194289 ship	n03417042 garbage truck
n02374451 horse	n03095699 container ship	n03173929 delivery truck
n02382948 racehorse	n02981792 catamaran	n04467665 trailer truck
n02379183 quarter horse	n03344393 fireboat	n03345487 fire engine

port the performance on contrasting covariate-shifted ID with training ID, although covariate-shifted ID are not OOD samples. We use three metrics to evaluate the OOD detection performance, which are detailed as follows: **1) FPR95** stands for false positive rate measured when true positive rate (TPR) sits at 95%. Intuitively, FPR95 measures the portion of samples that are falsely recognized as in-distribution data when most true in-distribution samples are recalled. **2) AUROC** refers to the Area Under the Receiver Operating Characteristic curve, which is concerned with both FPR and TPR. **3) AUPR** means the Area Under the Precision-Recall curve, which considers both precision and recall. For FPR95, the lower the value, the better the model. For AUROC and AUPR, the higher the value, the better the model.

## 5 Experiments

**Implementation Details** We conduct experiments on the three proposed FS-OOD benchmarks, *i.e.*, DIGITS, OBJECTS, and COVID. In terms of architectures, we use LeNet-5 [39] for DIGITS and ResNet-18 [40] for both OBJECTS and COVID. All models are trained by the SGD optimizer with a weight decay of  $5 \times 10^{-4}$  and a momentum of 0.9. For DIGITS and OBJECTS, we set the initial learning rate to 0.1, which is decayed by the cosine annealing rule, and the total epochs to 100. For COVID benchmark, the initial learning rate is set to 0.001 and the model is trained for 200 epochs. When fine-tuning for source-awareness

enhancement, the learning rate is set to 0.005 and the total number of epochs is 10. The batch size is set to 128 for all benchmarks.

Notice that the baseline implementations of ODIN [3] and MDS [5] require validation set for hyperparameter tuning, we spare a certain portion of near-OOD for validation. More specifically, we use 1,000 notMNIST images for the DIGITS benchmark, 1,000 CIFAR-100 images for the OBJECTS benchmark, and 54 images from CT-SCAN dataset for the COVID benchmark. The proposed method SEM relies on the hyperparameter of  $M = 3$  for low-layer  $p(x_n)$  and number of classes for high-layer  $p(x)$  in Gaussian mixture model. For output features with dimensions over 50, PCA is performed to reduce the dimensions to 50.

### 5.1 Results on FS-OOD Setting

We first discuss the results on near- and far-OOD datasets. Table 2 summarizes the results where the proposed SEM is compared with current state-of-the-art methods including MSP [2], ODIN [3], Mahalanobis distance score (MDS), and Energy-based OOD [4].

**DIGITS Benchmark** For the DIGITS benchmark, SEM gains significant improvements in all metrics (FPR95, AUROC, and AUPR). A huge gain is observed on notMNIST, which is a challenging dataset due to its closeness in background to the training ID MNIST. While none of the previous softmax/logits-based methods (*e.g.*, MSP, ODIN, and EBO) are capable to solve the notMNIST problem, the proposed SEM largely reduces the FPR95 metric from 99% to 10.93%, and the AUROC is increased from around 30% to beyond 95%. One explanation of the clear advantage is that, the previous output-based OOD detection methods largely depend on the covariate shift to detect OOD samples, while the feature-based MDS (partly rely on top-layer semantic-aware features) and the proposed SEM uses more semantic information, which is critical to distinguish MNIST and notMNIST. In other words, in the MNIST/notMNIST scenario where ID and OOD have high visual similarity, large dependency on covariate shift while ignorance on the semantic information will lead to the failure of OOD separation. Similar advantages are also achieved with the other near-OOD dataset.

**OBJECTS Benchmark** Similar to DIGITS benchmark, the proposed SEM surpasses the previous state-of-the-art methods on the near-OOD scenario of the OBJECTS benchmark, especially on the more robust metrics of AUROC and AUPR. However, the performance gap is not as large as DIGITS. One explanation

Table 2: Comparison between previous state-of-the-art methods and the proposed SEM score on FS-OOD benchmarks. The proposed SEM obtains a consistently better performance on most of the metrics than MSP [2], ODIN [3], Energy-based OOD (EBO) score [4], and Mahalanobis Distance Score (MDS) [5], especially on the near-OOD scenarios.

	FPR95 ↓					AUROC ↑					AUPR ↑				
	MSP	ODIN	EBO	MDS	SEM	MSP	ODIN	EBO	MDS	SEM	MSP	ODIN	EBO	MDS	SEM
<b>- DIGITS</b> (Training ID: <u>MNIST</u> , Covariate-Shifted ID: <u>USPS &amp; SVHN</u> )															
notmnist	99.97	99.95	99.99	78.83	<b>10.93</b>	32.54	29.04	25.49	79.10	<b>96.74</b>	67.33	65.97	63.97	90.60	<b>98.54</b>
FashionMNIST	99.90	99.97	99.98	94.68	<b>68.63</b>	39.71	38.51	37.64	60.42	<b>80.20</b>	82.40	82.16	81.57	88.84	<b>94.38</b>
Mean (Near-OOD)	99.93	99.96	99.98	86.75	<b>39.78</b>	36.12	33.77	31.56	69.76	<b>88.47</b>	74.87	74.06	72.77	89.72	<b>96.46</b>
Texture	94.89	94.65	98.40	<b>87.46</b>	90.90	64.34	64.02	65.02	72.42	<b>74.45</b>	94.40	94.32	94.47	95.81	<b>96.12</b>
CIFAR-10	98.01	98.38	99.62	95.47	<b>91.57</b>	52.22	51.15	50.95	67.96	<b>69.29</b>	87.26	86.84	86.36	91.74	<b>92.06</b>
Tiny-ImageNet	97.98	98.23	99.58	96.20	<b>93.39</b>	52.94	51.98	51.89	64.31	<b>67.54</b>	87.51	87.15	86.72	90.71	<b>91.58</b>
Places365	98.68	98.78	99.65	98.06	<b>94.15</b>	50.22	49.30	48.95	65.42	<b>67.63</b>	67.11	66.51	65.41	76.64	<b>77.61</b>
Mean (Far-OOD)	97.39	97.51	99.31	94.30	<b>92.50</b>	54.93	54.11	54.20	67.53	<b>69.73</b>	84.07	83.71	83.24	88.73	<b>89.34</b>
<b>- OBJECTS</b> (Training ID: <u>CIFAR-10</u> , Covariate-Shifted ID: <u>CIFAR-10-C &amp; ImageNet-10</u> )															
CIFAR-100	89.44	87.51	<b>83.84</b>	86.28	86.96	70.17	60.29	63.85	72.05	<b>74.70</b>	88.28	81.11	83.51	89.42	<b>90.64</b>
Tiny-ImageNet	88.22	88.13	<b>81.58</b>	87.45	86.59	72.92	62.07	67.97	72.94	<b>76.76</b>	90.04	82.49	86.30	89.96	<b>91.86</b>
Mean (Near-OOD)	88.83	87.82	<b>82.71</b>	86.87	86.77	71.55	61.18	65.91	72.50	<b>75.73</b>	89.16	81.80	84.91	89.69	<b>91.25</b>
MNIST	93.54	<b>82.04</b>	92.23	84.59	99.70	66.98	70.31	54.55	<b>77.04</b>	75.69	52.66	49.58	34.14	65.31	<b>76.61</b>
FashionMNIST	88.08	<b>68.73</b>	72.40	77.17	93.72	73.78	<b>80.98</b>	76.50	80.33	79.40	90.15	91.53	89.80	92.28	<b>93.14</b>
Texture	85.64	<b>72.91</b>	75.57	72.98	82.15	74.18	70.14	68.63	72.02	<b>79.69</b>	93.34	89.97	89.51	88.46	<b>95.48</b>
CIFAR-100-C	87.26	84.26	<b>83.64</b>	85.53	83.92	74.12	67.51	68.37	68.13	<b>78.89</b>	89.74	83.97	85.54	82.97	<b>92.07</b>
Mean (Far-OOD)	88.63	76.98	80.96	<b>80.07</b>	89.87	72.27	72.23	67.01	74.38	<b>78.42</b>	81.47	78.76	74.75	82.25	<b>89.33</b>
<b>- COVID</b> (Training ID: <u>BIMCV</u> , Covariate-Shifted ID: <u>ActMed &amp; Hannover</u> )															
CT-SCAN	99.80	93.06	97.35	99.39	<b>2.24</b>	11.31	26.57	13.14	81.21	<b>99.51</b>	52.92	57.44	53.34	94.31	<b>99.80</b>
XRrayBone	97.00	55.50	42.00	100.00	<b>14.50</b>	32.08	64.73	77.80	78.72	<b>94.97</b>	76.95	86.11	91.68	96.67	<b>98.95</b>
Mean (Near-OOD)	98.40	74.28	69.67	99.69	<b>8.37</b>	21.70	45.65	45.47	79.96	<b>97.24</b>	64.94	71.77	72.51	95.49	<b>99.37</b>
MNIST	98.30	65.14	0.35	100.00	<b>0.00</b>	24.89	65.37	99.91	80.81	<b>100.00</b>	1.07	2.33	95.90	81.11	<b>100.00</b>
CIFAR-10	96.32	84.61	94.67	98.02	<b>85.58</b>	41.12	57.70	45.23	77.05	<b>52.50</b>	8.73	12.17	9.77	61.14	<b>11.27</b>
Texture	98.39	94.59	87.06	56.38	<b>27.57</b>	22.63	31.13	34.95	89.84	<b>90.94</b>	11.43	12.59	13.14	85.50	<b>64.71</b>
Tiny-ImageNet	97.78	90.26	92.73	92.11	<b>44.99</b>	30.26	42.76	32.69	81.99	<b>83.42</b>	7.42	8.90	7.65	77.94	<b>31.19</b>
Mean (Far-OOD)	97.70	83.65	68.70	86.63	<b>39.54</b>	29.73	49.24	53.20	82.42	<b>81.72</b>	7.16	9.00	31.62	76.42	<b>51.79</b>

is that images in OBJECTS benchmark are more complex than DIGITS, leading the neural networks to be more semantics-orientated. Therefore, more semantic information is encoded in the previous output-based methods. Nevertheless, the proposed SEM method still outperforms others on most of the metrics. We also notice that SEM score does not reach the best performance on MNIST and FashionMNIST. One explanation is that two black-and-white images in these two datasets inherently contain significant covariate shifts comparing to both training ID and covariate-shifted ID, so that the scores that efficient on covariate shift detection (*e.g.*, ODIN) can also achieve good results on these datasets. However, these methods fail in near-OOD scenario, as they might believe CIFAR-10-C should be more likely to be OOD than CIFAR-100.

**COVID Benchmark** In this new and real-world application of OOD detection, the proposed SEM score achieves an extraordinary performance on all metrics, which surpasses the previous state-of-the-art methods by a large margin in both near and far-OOD scenarios. The result also indicates that previous output-based methods generally breaks down on this setting, *e.g.*, their FPR@95 scores are generally beyond 90% in near-

OOD setting which means ID and OOD are totally mixed. However, the proposed SEM achieves around 10% in near-OOD setting. On far-OOD samples, the output-based methods are still unable to be sensitive to the ID/OOD discrepancy. The phenomenon matches the performance in DIGITS dataset, where the training data is simple and the logits might learn much non-semantic knowledge to be cancelled out.

**Observation Summary** We summarize the following two take-away messages from the experiments on all three FS-OOD benchmarks: **1)** SEM score performs consistently well on near-OOD, which classic output-based methods (*e.g.*, MSP, ODIN, EBO) majorly fail on. The reason can be that output-based methods use too much covariate shift information for OOD detection, which by nature cannot distinguish between covariate-shifted ID and near-OOD. The proposed SEM score also outperforms the similar feature-based baseline MDS. **2)** SEM score sometimes underperforms on far-OOD, with a similar reason that classic OOD detectors use covariate shift to distinguish ID and OOD, which is sometimes sufficient to detect far-OOD samples. Nevertheless, SEM reaches more balanced good results on near-OOD and far-OOD.

Table 3: Comparison between previous state-of-the-art methods, the proposed SEM score, and the low-level probabilistic component  $p(\mathbf{x}_n)$  on classic OOD benchmarks, without the existence of covariate-shifted ID set. The previous methods of MSP [2], ODIN [3], EBO score [4], and MDS [5] reaches a good results on the classic benchmark. However, the value of  $p(\mathbf{x}_n)$  can exceed all the previous methods and achieve a near-perfect result across all the metrics, showing that only taking covariate shift score can completely solve the classic OOD detection benchmark, which, in fact, contradicts the goal of OOD detection. This phenomenon also advocates the significance of the proposed FS-OOD benchmark.

	FPR95 ↓						AUROC ↑						AUPR ↑					
	MSP	ODIN	MDS	EBO	SEM	$p(\mathbf{x}_n)$	MSP	ODIN	MDS	EBO	SEM	$p(\mathbf{x}_n)$	MSP	ODIN	MDS	EBO	SEM	$p(\mathbf{x}_n)$
<b>- DIGITS (ID: MNIST)</b>																		
notMNIST	43.09	37.70	44.06	1.77	2.64	<b>0.78</b>	88.77	89.85	88.44	99.67	99.50	<b>99.79</b>	75.72	77.83	75.97	99.36	99.09	<b>99.57</b>
FashionMNIST	2.54	1.08	1.05	0.27	40.09	<b>0.00</b>	99.44	99.70	99.72	99.90	95.02	<b>99.94</b>	99.64	99.77	99.76	99.94	97.63	<b>99.97</b>
Mean (Near-OOD)	20.05	13.48	20.54	2.68	27.85	<b>0.46</b>	96.06	96.97	95.85	99.49	93.85	<b>99.78</b>	94.07	94.72	92.66	99.40	93.23	<b>99.73</b>
Texture	2.43	0.94	0.67	0.23	90.69	<b>0.02</b>	99.34	99.75	99.81	99.93	77.26	<b>99.91</b>	99.58	99.84	99.84	99.96	87.56	<b>99.95</b>
CIFAR-10	7.05	3.06	3.18	0.18	54.43	<b>0.00</b>	98.68	99.31	99.30	99.88	94.19	<b>99.97</b>	98.72	99.27	99.12	99.88	95.86	<b>99.97</b>
Tiny-ImageNet	6.28	2.93	3.13	0.55	59.52	<b>0.00</b>	98.78	99.36	99.37	99.79	93.70	<b>99.96</b>	98.78	99.33	99.25	99.79	95.54	<b>99.96</b>
Places365	9.92	4.59	4.12	0.45	58.07	<b>0.00</b>	98.19	99.06	99.17	99.81	93.82	<b>99.96</b>	94.87	97.01	96.84	99.42	91.32	<b>99.88</b>
Mean (Far-OOD)	6.45	2.92	2.87	0.36	53.03	<b>0.00</b>	98.77	99.36	99.39	99.84	94.18	<b>99.96</b>	98.00	98.84	98.74	99.76	95.09	<b>99.94</b>

## 5.2 Results on Classic OOD Detection Setting

Table 3 shows the performance on the classic OOD detection benchmark. The result shows that without the introduction of covariate-shifted ID data, the previous methods reach a near-perfect performance on the classic benchmark, which matches the reported results in their origin papers. However, by comparing with Table 2, their performance significantly breakdown when covariate-shifted ID is introduced, showing the fragility of previous methods, and therefore we advocate the more realistic FS-OOD benchmark. Furthermore, we also report the results that by using the value of  $p(\mathbf{x}_n)$ , the score from low-layer feature statistics for detecting covariate shift is shown surprisingly effective on classic OOD benchmark, which exceeds all the previous methods and achieve a near-perfect result across all the metrics. This phenomenon shows that only taking covariate shift score can completely solve the classic OOD detection benchmark with MNIST, which, in fact, contradicts the goal of OOD detection. It also advocates the significance of the proposed FS-OOD benchmark.

## 5.3 Ablation Study

In this section, we validate the effectiveness of the main components that contribute to the proposed SEM score, and also analyze the effects of fine-tuning scheme for source-awareness enhancement. All the experiments in this part are conducted on the DIGITS benchmark.

**Components of SEM** According to Equation 2 in the Section 3, SEM score can be decomposed by the estimations of  $p(\mathbf{x})$  and  $p(\mathbf{x}_n)$ . While our final SEM score uses output flattened features of the CNN model for  $p(\mathbf{x})$  estimation and low-layer feature statistics for

Table 4: Ablation study on the SEM components. AUROC is reported for performance evaluation. Several options can be applied to estimate  $p(\mathbf{x}_n)$  and  $p(\mathbf{x})$  in Equation 3. FS denotes the usage of feature statistics, and FF denotes flattened features. T and L means top-/low-layer feature, *e.g.*, L-FS means low-layer feature statistics. The results show the effectiveness of our SEM score.

#	$p(\mathbf{x}_n)$			$p(\mathbf{x})$	NearOOD	FarOOD
	T-FS	L-FF	L-FS	T-FF		
1				✓	87.28	60.80
2	✓				87.28	60.80
3	✓			✓	-	-
4		✓			51.81	51.81
5		✓		✓	86.54	61.26
6			✓		70.27	72.58
7			✓	✓	<b>88.47</b>	<b>69.73</b>

$p(\mathbf{x}_n)$ , there are actually several options for the estimation, which is discussed in Table 4. In this analysis, we set top flattened features as the default usage for  $p(\mathbf{x})$  and only explore  $p(\mathbf{x}_n)$ , which is the key part of SEM score.

Exp#1 shows the result that only uses  $p(\mathbf{x})$  as the final score, which can be interpreted as a simple method using GMM to estimate ID likelihood on the final-layer features. Compared to the MDS result in Table 2, this simple method already obtains a better performance on near-OOD. Notice that we use LeNet-5 on DIGITS, the final-layer features are identical to their feature statistics (ref. Exp#2). Therefore, everything is cancelled out if  $p(\mathbf{x}_n)$  is top-layer feature statistics (ref. Exp#3).

Exp#4 and Exp#6 shows comparison between using low-layer flattened features (L-FF) and low-layer feature statistics (L-FS) only. The performance on detecting covariate-shifted ID shows that both L-FF and L-FS

Table 5: **Ablation study on the fine-tuning scheme for source-awareness enhancement.** AUROC is reported for performance evaluation. #1 reports the performance before fine-tuning.  $\mathcal{L}_{src}(x)$  means fine-tuning without negative augmented data.  $\mathcal{L}_{src}(x')$  means only data with negative augmentation is used. The results show the effectiveness of each training loss.

#	$\mathcal{L}_{cls}$	$\mathcal{L}_{src}(x)$	$\mathcal{L}_{src}(x')$	NearOOD	FarOOD
1				83.03	56.65
2	✓			86.55	64.61
3	✓	✓		87.42	68.40
4	✓		✓	87.27	67.92
5	✓	✓	✓	<b>88.47</b>	<b>69.73</b>

Table 6: **Hyperparameter Selection of the Number of GMM Components  $K$ .** The result shows that  $M = 3$  in low-layer statistics and  $M = 10$  for top-layer features (equal to number of classes) can reach the best results in MNIST benchmark.

#	$p(x_n)$		$p(x)$		NearOOD	FarOOD
	M=1	M=3	M=10	M=20		
1	✓			✓	86.24	64.94
2	✓		✓		85.81	60.17
3		✓		✓	84.47	63.61
4		✓	✓		<b>88.47</b>	<b>69.73</b>

have significant sensitivity to covariate shifts, but with a poor performance on FS-OOD detection. The result indicates that with only the usage of low-level features, the score has a strong correlation to covariate shift but barely contains semantic information, and the feature statistics show the stronger characteristics compared to flattened feature. This observation indicates our selection of low-level feature statistics for estimating  $p(x_n)$ , which is further supported by the results of Exp#5 and Exp#7, and visually illustrated by Figure 4.

**Fine-Tuning Scheme** Here we evaluate the designed fine-tuning scheme of SEM. As elaborated in Section 3.2, this learning procedure is designed to enhance the source-aware compactness. Specifically, a source-awareness enhancement loss  $\mathcal{L}_{src}$  is proposed to aggregate the ID training data and separate from the generated negative augmented images at the same time. Table 5 demonstrates the effectiveness of the fine-tuning scheme. When combining both in-distribution training and negative augmented data training, our framework achieves the best performance.

**Hyperparameter of  $M$**  Table 6 shows the analysis of hyperparameter  $M$ . In the DIGITS dataset,  $M = 3$  leads to a slightly better performance comparing to other choices. Nevertheless, the overall difference among

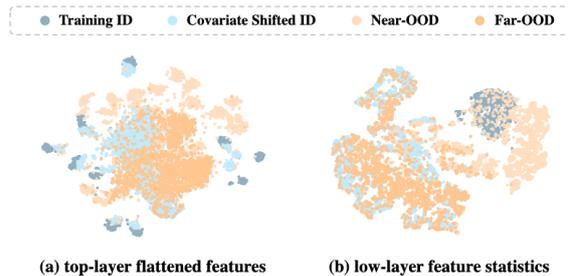


Fig. 4: T-SNE visualization on DIGITS. It suggests that low-layer feature statistics capture non-semantic information, and top-layer features capture both semantic and non-semantic information.

various  $M$  is not obvious on near-OOD, showing that the model is robust to the hyperparameter.

## 6 Discussion and Conclusion

Existing OOD detection literature has shown mostly relied on covariate shift even though they are intended to detect semantic shift. This is very effective when test OOD data only come from the far-OOD group—where the covariate shift is large and is further exacerbated by semantic shift, so using covariate shift as a measure to detect OOD fares well. However, when it comes to near-OOD data, especially with covariate-shifted ID (*i.e.*, data experiencing covariate shift but still belonging to the same in-distribution data), current state-of-the-art methods would suffer a significant drop in performance, as shown in the experiments.

We find the gap is caused by a shortcoming in existing evaluation benchmarks: they either exclude covariate-shifted data during testing or treat them as OOD, which is conceptually contradictory with the primary goal that a machine learning model should generalize beyond the training distribution. To fill the gap, we introduce a new problem setting that better matches the design principles of machine learning models: they should be robust in terms of good generalization to covariate-shifted datasets, and trustworthy as they also need to be capable of detecting abnormal semantic shift.

The empirical results suggest that current state-of-the-art methods rely too heavily on covariate shift and hence could easily mis-classify covariate-shifted ID data as OOD data. In contrast, our SEM score function, despite having a simple design, provides a more reliable measure for solving full-spectrum OOD detection.

In fact, to detecting samples with covariate shift, we find that a simple probabilistic model using low-level feature statistics can reach a near-perfect result.

**Outlook** As the OOD detection community getting common awareness of the saturated performance

problem of classic OOD benchmarks, several works have taken one-step further towards the more realistic setting and proposed large-scale benchmarks [41, 42]. However, this paper shows that even under the classic MNIST/CIFAR-scale OOD benchmarks, current OOD methods in fact cannot achieve satisfactory results when the generalization ability is required. We hope that the future OOD detection works could also consider the generalization capability on covariate-shifted ID data, in parallel to exploring larger-scale models and datasets.

**Broader Impacts** Our research aims to improve the robustness of machine learning systems in terms of the capability to safely handle abnormal data to avoid catastrophic failures. This could have positive impacts on a number of applications, ranging from consumer (e.g., AI-powered mobile phones) to transportation (e.g., autonomous driving) to medical care (e.g., abnormality detection). The new problem setting introduced in the paper includes an important but largely missing element in existing research, namely data experiencing covariate shift but belonging to the same in-distribution classes. We hope the new setting, along with the simple approach based on SEM and the findings presented in the paper, can pave the way for future research for more reliable and practical OOD detection.

## References

- Jingkang Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu. Generalized out-of-distribution detection: A survey. *arXiv preprint arXiv:2110.11334*, 2021.
- Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2017.
- Shiyu Liang, Yixuan Li, and Rayadurgam Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2017.
- Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- Jie Ren, Peter J Liu, Emily Fertig, Jasper Snoek, Ryan Poplin, Mark Depristo, Joshua Dillon, and Balaji Lakshminarayanan. Likelihood ratios for out-of-distribution detection. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- Joost Van Amersfoort, Lewis Smith, Yee Whye Teh, and Yarin Gal. Uncertainty estimation using a single deep deterministic neural network. In *ICML*, 2020.
- Chandramouli Shama Sastry and Sageev Oore. Detecting out-of-distribution examples with gram matrices. In *ICML*, 2020.
- Kaiyang Zhou, Ziwei Liu, Yu Qiao, Tao Xiang, and Chen Change Loy. Domain generalization in vision: A survey. *arXiv preprint arXiv:2103.02503*, 2021.
- Yifei Ming, Hang Yin, and Yixuan Li. On the impact of spurious correlation for out-of-distribution detection. *arXiv preprint arXiv:2109.05642*, 2021.
- Hyunsun Choi, Eric Jang, and Alexander A Alemi. WAIC, but why? Generative ensembles for robust anomaly detection. *arXiv preprint arXiv:1810.01392*, 2018.
- Apoorv Vyas, Nataraj Jammalamadaka, Xia Zhu, Dipankar Das, Bharat Kaul, and Theodore L Willke. Out-of-distribution detection using an ensemble of self supervised leave-out classifiers. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- Joan Serra, David Álvarez, Vicenç Gómez, Olga Sliozovskaia, José F Núñez, and Jordi Luque. Input complexity and out-of-distribution detection with likelihood-based generative models. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2020.
- Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. Deep anomaly detection with outlier exposure. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2019.
- Yi Li and Nuno Vasconcelos. Background data resampling for outlier-aware classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- Jingkang Yang, Haoqi Wang, Litong Feng, Xiaopeng Yan, Huabin Zheng, Wayne Zhang, and Ziwei Liu. Semantically coherent out-of-distribution detection. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021.
- Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2018.
- Chuang Lin, Zehuan Yuan, Sicheng Zhao, Peize Sun, Changhu Wang, and Jianfei Cai. Domain-invariant disentangled network for generalizable object detection. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021.
- Xingchao Peng, Zijun Huang, Ximeng Sun, and Kate Saenko. Domain agnostic learning with disentangled representations. In *Proceedings of International Conference on Machine Learning (ICML)*, 2019.
- Kaiyang Zhou, Yongxin Yang, Yu Qiao, and Tao Xiang. Domain generalization with mixstyle. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2021.
- Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2013.
- Abhishek Sinha, Kumar Ayush, Jiaming Song, Burak Uzkent, Hongxia Jin, and Stefano Ermon. Negative data augmentation. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2021.
- Yann LeCun, Corinna Cortes, and Christopher JC Burges. The mnist database of handwritten digits, 1998. <http://yann.lecun.com/exdb/mnist>, 1998.
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bisaccho, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. In

- Proceedings of NIPS Workshop on Deep Learning and Unsupervised Feature Learning*, 2011.
26. Jonathan J. Hull. A database for handwritten text recognition research. *IEEE Transactions on pattern analysis and machine intelligence*, 16(5):550–554, 1994.
  27. Yaroslav Bulatov. NotMNIST dataset. <http://yaroslavvb.blogspot.com/2011/09/notmnist-dataset.html>, 2011.
  28. Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
  29. 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 (CVPR)*, 2014.
  30. Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. *CiteSeer*, 2009.
  31. Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large-scale visual recognition challenge. *International Journal of Computer Vision (IJCV)*, 2015.
  32. Bolei Zhou, Agata Lapedrizza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2017.
  33. Beatriz Garcia Santa Cruz, Matías Nicolás Bossa, Jan Sölter, and Andreas Dominik Husch. Public covid-19 x-ray datasets and their impact on model bias—a systematic review of a significant problem. *Medical image analysis*, 2021.
  34. Maria de la Iglesia Vayá, Jose Manuel Saborit, Joaquim Angel Montell, Antonio Pertusa, Aurelia Bustos, Miguel Cazorla, Joaquin Galant, Xavier Barber, Domingo Orozco-Beltrán, Francisco García-García, et al. Bimcv covid-19+: a large annotated dataset of rx and ct images from covid-19 patients. *arXiv preprint arXiv:2006.01174*, 2020.
  35. Linda Wang, Zhong Qiu Lin, and Alexander Wong. Covid-net: a tailored deep convolutional neural network design for detection of covid-19 cases from chest x-ray images. *Scientific Reports*, 2020.
  36. Hinrich B Winther, Hans Laser, Svetlana Gerbel, Sabine K Maschke, Jan B Hinrichs, Jens Vogel-Claussen, Frank K Wacker, Marius M Höper, and Bernhard C Meyer. Covid-19 image repository. 2020. URL [https://figshare.com/articles/COVID-19\\_Image\\_Repository/12275009](https://figshare.com/articles/COVID-19_Image_Repository/12275009).
  37. RSNA. RSNA Pediatric Bone Age Challenge (2017), 2017.
  38. Xingyi Yang, Xuehai He, Jinyu Zhao, Yichen Zhang, Shanghang Zhang, and Pengtao Xie. Covid-ct-dataset: a ct scan dataset about covid-19. *arXiv preprint arXiv:2003.13865*, 2020.
  39. Yann LeCun et al. Lenet-5, convolutional neural networks. <http://yann.lecun.com/exdb/lenet>, 2015.
  40. Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
  41. Haoqi Wang, Zhizhong Li, Litong Feng, and Wayne Zhang. Vim: Out-of-distribution with virtual-logit matching. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
  42. Anugya Srivastava, Shriya Jain, and Mugdha Thigle. Out of distribution detection on imagenet-o. *arXiv preprint arXiv:2201.09352*, 2022.