
IGOOD: An Information Geometry Approach to Out-of-Distribution Detection

Eduardo D. C. Gomes, Florence Alberge, Pierre Duhamel and Pablo Piantanida

Laboratoire des signaux et systèmes (L2S)
Université Paris-Saclay CNRS CentraleSupélec
91190, Gif-sur-Yvette, France.

Abstract

Reliable out-of-distribution (OOD) detection is a fundamental step towards a safer implementation of modern machine learning (ML) systems under distribution shift. In this paper, we introduce IGOOD, an effective method for detecting OOD samples. IGOOD applies to any pre-trained neural network, works under different degrees of access to the ML model, does not require OOD samples or assumptions on the OOD data, but can also benefit (if available) from OOD samples. By building on the geodesic (Fisher-Rao) distance between the underlying data distributions, our discriminator combines confidence scores from the logits outputs and the learned features of a deep neural network. Empirically, we show that IGOOD is competitive and often outperforms state-of-the-art methods by a large margin on a variety of networks architectures and datasets.

1 Introduction

Out-Of-Distribution (OOD) or novelty detection is one of the main objectives in conceiving reliable ML systems [1]. An application of these methods arises in monitoring ML-based online services for drifting distributions. Tracking changes in the underlying data distribution is challenging as they contain unusual (irregular or unexpected) events and have large dimensions. For instance, the detection of such anomalous observations will have to rely on the intrinsic properties of the ML models, methods, and algorithms based on their statistical behavior in the presence of in-distribution data. Classic approaches to OOD detection consist of deriving metrics for detecting those abnormalities from the lens of ML models (e.g., softmax output, latent representations across layers), provided that the data is high dimensional and often only a single test example is available.

1.1 Contributions

In this paper, we propose IGOOD, a new unified and effective method to perform OOD detection by rigorously exploring the information-geometric properties of the feature space on various depths of a DNN. IGOOD provides a flexible framework that applies to any pre-trained softmax neural classifier. A key ingredient of IGOOD is the Fisher-Rao distance. This distance is used as an effective differential geometry tool for clustering in the context of multivariate Gaussian pdfs [30, 37]. In our context, we measure the dissimilarity between probability (in and out) distributions, as the length of the shortest path within the manifold induced by the underlying class of distributions (i.e., the softmax probabilities of the neural classifier or the densities modeling the learned representations across the layers). By doing so, we can explore statistical invariances of the geometric properties of the learned features [5]. Our method adapts to various scenarios depending on the level of information access of the DNN, uses only in-distribution samples but can also benefit (if available) of OOD samples.

1.2 Related works

OOD detectors are binary classifiers that discriminate in- and out-of-distribution samples. A few works [36, 14, 4, 26, 40, 39, 12] propose retraining the base (or an auxiliary) model with synthetic or ground truth OOD samples to serve as a classifier and as an OOD discriminator. Disposing of both OOD and in-distribution samples during training enables the hidden layers to learn representations to facilitate OOD detection. These methods will not be compared to ours in this work, as they entail retraining or modifying the base neural network by using OOD data to further train parameters. Moreover, this assumes that OOD samples are stationary, which is an unrealistic assumption in practical scenarios. The work [27] demonstrates failure modes of OOD detection methods to better understand how to improve them, especially how spurious features like image background can vastly degrade detection performance. references [17, 6, 38, 42, 33, 45, 25] study OOD detection in the context of generative models. Open set recognition [19], outlier or anomaly detection [29], concept drift detection [31], and adversarial attacks detection [10, 23] are related topics.

WHITE-BOX scenario. This class of OOD detectors allows discriminators to have access to all intermediate layer outputs. Naturally, they have access to more information than **BLACK-BOX** or **GREY-BOX** techniques, which provide detection based only on the network's outputs, i.g., MSP [15], ODIN [21], and the free-energy [22] based methods to name a few. reference [34] proposes high order Gram matrices to perform OOD detection, by computing class-conditional pairwise feature correlations across the hidden layers of the network. The work [20] models the latent features' outputs of DNN models as a class-conditional Gaussian mixture distribution with tied covariance matrix and class-conditional mean vectors. They calculate the Mahalanobis distance between an OOD sample as a single estimator of the mean of a class-conditional Gaussian distribution with covariance matrix estimated on the entire training set. The importance of each feature component and hyperparameters are tuned using validation data. The work [32] modifies the Mahalanobis distance-based OOD detector [20] to improve near-OOD detection by reducing the importance of features shared by in- and out-of-distribution data.

2 IGEOD: OOD Detection using the Fisher-Rao Distance

In this section, we introduce IGEOD, a flexible framework for OOD detection. IGEOD is implemented in two ways: at the level of the logits using temperature scaling (Section 2.1), and layer-wise level (Section 2.2). The key ingredient of IGEOD is the Fisher-Rao (F-R) distance [2]. This distance measures the dissimilarity between two probability models within a class of probability distributions by calculating the geodesic distance between two points on the learned manifold.

2.1 IGEOD score using the softmax probability

For the classification problem, we can take the temperature scaled softmax function (Eq.1) as an approximation of a class-conditional probability distribution:

$$q(y|f(x); T) = \frac{\exp(f_y(x)/T)}{\sum_{y \in \mathcal{Y}} \exp(f_y(x)/T)}; \quad (1)$$

where $f: X \rightarrow \mathbb{R}^C$ is a vectorial function with f_1, f_2, \dots, f_C and $f_y(\cdot)$ denotes the y -th logits output value of the DNN classifier. The F-R distance d_{FR}^{Logits} between two softmax probability distributions can be shown by

$$d_{FR}^{Logits}(q(f(x)); q(f(x^0))) = 2 \arccos \sqrt{\sum_{y \in \mathcal{Y}} q_y(f(x)) q_y(f(x^0))}; \quad (2)$$

Class conditional centroid estimation. We model the training dataset class-conditional posterior distribution by calculating the centroid of the logits representations of this set. Precisely, we compute the centroid for the logits of each class of the in-distribution training dataset \mathcal{D}_N corresponding to

¹We refer the reader to the appendix (see Section A).

the F-R distance, i.e.,

$$d_{FR}^{(y)} = \min_{\theta \in \mathcal{C}} \frac{1}{N_y} \sum_{x_i: y_i=y} d_{FR}^{Logits}(\theta; q(j|x_i); q(j|y)); \quad (3)$$

where N_y is the amount of training examples with label y . We optimize this expression using SGD algorithm, where the parameter to be tuned is in the logits space.

OOD detection score. We propose the F-R distance-based OOD detection score ($FR_0(x)$) on the space of the logits to be the sum of the distances between x and each individual class conditional centroid y calculated by Eq. (3). We denote it as follows:

$$FR_0(x) = \sum_{y \in \mathcal{Y}} d_{FR}^{Logits}(\theta; q(j|x); q(j|y)); \quad (4)$$

We obtained better performance by taking the sum instead of the minimal distance. A likely explanation for this would be that Eq. (4) leverages useful information related to the example's confidence score for each class.

2.2 IGEOD score leveraging latent features

For each layer, we define a set of class-conditional Gaussian distributions with diagonal standard deviation matrix $\Sigma_j^{(l)}$ and class-conditional mean $\mu_j^{(l)}$, where $l \in \{1, \dots, C\}$ and j is the index of the latent feature. We compute the empirical estimates of these parameters according to

$$\mu_j^{(l)} = \frac{1}{N_y} \sum_{x_i: y_i=y} f_j^{(l)}(x_i); \quad \text{and} \quad \Sigma_j^{(l)} = \text{diag} \left(\frac{1}{N_y} \sum_{x_i: y_i=y} f_j^{(l)}(x_i)^2 - \mu_j^{(l)2} \right); \quad (5)$$

where $j \in \{1, \dots, k\}$, k is the size of feature, and $f_j^{(l)}(\cdot)$ is the output of the network for feature j . The F-R distance d_{FR}^{Gauss} between two univariate Gaussian pdfs $\mathcal{N}(\mu_1; \Sigma_1)$ and $\mathcal{N}(\mu_2; \Sigma_2)$ is given by²

$$d_{FR}(\mathcal{N}(\mu_1; \Sigma_1); \mathcal{N}(\mu_2; \Sigma_2)) = \frac{1}{2} \log \frac{\Sigma_1^{1/2} \Sigma_2^{1/2} + \Sigma_1^{1/2} \mu_2 + \Sigma_2^{1/2} \mu_1}{\Sigma_1^{1/2} \mu_1 + \Sigma_2^{1/2} \mu_2}; \quad (6)$$

where $(\mu_i; \Sigma_i)$ is a 2-dimensional vector with components μ_i and Σ_i and $\|\cdot\|_2$ is the 2-norm. Similarly, the F-R distance d_{FR}^{Gauss} between two multivariate Gaussian pdfs with diagonal standard deviation matrix is derived from the univariate case and is given by

$$d_{FR}^{Gauss}(\mathcal{N}(\mu; \Sigma); \mathcal{N}(\mu_0; \Sigma_0)) = \frac{1}{2} \sum_{i=1}^k d_{FR}(\mathcal{N}(\mu_i; \Sigma_{ii}); \mathcal{N}(\mu_{0i}; \Sigma_{0i})); \quad (7)$$

where k is the cardinality of the distribution $\mathcal{N}(\mu; \Sigma)$ and $\mathcal{N}(\mu_0; \Sigma_0)$, μ_i is the i -th component of the vector μ , and Σ_{ii} is the entry of index i of the standard deviation matrix.

Experimental support for a diagonal Gaussian mixture model. We observed that the latent features covariance matrices are often conditioned and are diagonal dominant. In other words, the condition number of the covariance matrix often diverges, and the magnitude of the diagonal entry in a row is greater than or equal to the sum of all the other entries in that row for most rows.

Fisher-Rao distance-based feature-wise confidence score. We derive a confidence score by applying the F-R distance between the test sample and the closest class-conditional diagonal Gaussian distribution. We can consider two situations: (i) we do not have access to any validation OOD data whatsoever. In this case, the natural choice is to model the test samples as Gaussian distribution with the same diagonal standard deviation as the learned representation, i.e.,

$$FR(x) = \min_{y \in \mathcal{Y}} d_{FR}^{Gauss}(x; \mu_j^{(l)}; \Sigma_j^{(l)}); \quad (8)$$

and (ii), we dispose of a validation OOD dataset on which the features' diagonal standard deviation matrices $\alpha_j^{(l)}$ and $\alpha_j^{(l)}$ can be estimated, as well as the quantity

$$FR^0(x) = \min_{y \in \mathcal{Y}} d_{FR}^{Gauss}(x; \mu_j^{(l)}; \alpha_j^{(l)}); \quad (9)$$

²The reader can be referred to the appendix (Section A) for the derivation of this distance.

Feature ensemble. Similarly to [20], we combine the confidence scores of the logits and low-level features through a linear combination to emphasize features that demonstrate a greater capacity for detecting abnormal samples. The following equation summarizes the score:

$$FR(x) = \alpha FR_0(x) + \beta FR(x) + \gamma FR^0(x) \quad (10)$$

3 Experimental Results

The experimental setup follows the setting established by [21] and [20]. We use two pre-trained deep neural networks architectures for image classification tasks: DenseNet-B101 and a ResNet-34 [11]. We take in-distribution data images from CIFAR-10 [8], CIFAR-100 and SVHN [28] datasets. For out-of-distribution data we use natural image examples from the datasets: Tiny-ImageNet [9], LSUN [44], Describable Textures Dataset [7], Chars74K [8], Places365 [6], iSUN [43] and a synthetic dataset generated from Gaussian noise. For models pre-trained on CIFAR-10, data from CIFAR-100 and SVHN are also considered OOD; for models pre-trained on CIFAR-100, data from CIFAR-10 and SVHN are considered OOD, and for models pre-trained on SVHN, the CIFAR-10 and CIFAR-100 datasets are considered OOD. Even though we ran experiments with image data, \mathcal{G}_{OOD} could be applied to any neural-based classification task.

We consider two tuning scenarios: one with data from adversarially generated (FGSM) samples from the training dataset, and another with data from the OOD test set. For the former, we tune hyperparameters for each method with generated data (pseudo OOD) and in-distribution data. While for the latter, we tune hyperparameters with 1,000 OOD data samples and in-distribution data. We derive two methods: \mathcal{G}_{OOD} , which is given by Eq(10) and considers that we have an estimate on the diagonal covariance matrix and mean vector from OOD data as additional information; and \mathcal{I}_{OOD} , which doesn't consider any prior on OOD data, i.e., $\beta = 0$ on Eq. (10).

Comparison with Mahalanobis. For each DNN model and in-distribution dataset pair, we report the average OOD detection performance for Mahalanobis [20], \mathcal{G}_{OOD} and \mathcal{I}_{OOD} . Table 1 validates the contributions of our techniques. We observe substantial performance improvement in all experiments for the left-hand side of the table, where \mathcal{G}_{OOD} outperforms Mahalanobis on average for all test cases, recording an improvement up to 23% on TNR at TPR-95%. To assess the consistency of \mathcal{G}_{OOD} to the choice of validation data, we measured the detection performance when hyperparameters are tuned only using in-distribution and generated adversarial data as observed in the right-hand side of Table 1. In this setup, \mathcal{G}_{OOD} improves by 2.5% the average TNR at TPR-95% across all datasets and models, but is sometimes outperformed by 2-3%.

Table 1: Average and standard deviation of OOD detection performance for White-Box settings. The abbreviation TNR-95%, C-10 and C-100 stands for TNR at TPR-95%, CIFAR-10 and CIFAR-100, respectively. The extended results can be found in Tables 6 and 7 in the appendix.

Model	In-dist.	Tuning on OOD data			Tuning on adversarial data		
		TNR-95% Mahalanobis	AUROC \mathcal{G}_{OOD} (ours)	AUROC	TNR-95% Mahalanobis	AUROC \mathcal{G}_{OOD} (ours)	AUROC
DenseNet	C-10	76.6 ₃₁ /92.6 ₁₄	92.1 ₁₂ /98.4 _{3.0}		75.9 ₃₀ /77.9 ₂₉	91.7 ₁₂ /94.0 _{9.0}	
	C-100	67.2 ₂₈ /90.2 ₂₁	90.2 ₁₃ /97.7 _{5.0}		60.4 ₃₄ /70.9 ₃₅	85.3 ₁₉ /90.8 ₁₃	
	SVHN	93.3 _{8.0} /98.0 _{2.0}	98.6 _{1.0} /99.6 _{0.1}		93.7 ₁₀ /92.2 _{9.0}	98.6 _{2.0} /98.4 _{1.0}	
ResNet	C-10	82.5 ₂₃ /91.6 ₁₆	96.5 _{4.0} /98.4 _{3.0}		78.6 ₂₄ /77.3 ₃₂	95.3 _{6.0} /90.0 ₁₅	
	C-100	70.4 ₃₀ /86.4 ₂₃	91.9 ₁₀ /97.1 _{5.0}		57.4 ₃₆ /65.1 ₃₃	86.9 ₁₃ /88.6 ₁₅	
	SVHN	96.8 _{6.0} /98.9 _{2.0}	99.2 _{1.0} /99.7 _{0.1}		96.3 _{8.0} /93.6 ₁₄	99.1 _{1.0} /98.4 _{3.0}	
Average and Std.		81.1 ₁₁ /92.9 _{4.0}	94.8 _{4.0} /98.5 _{1.0}		77.0 ₁₅ /79.5 ₁₀	92.8 _{5.4} /93.4 _{3.9}	

Ablation study. The \mathcal{I}_{OOD} score has three components: FR_0 , FR , and FR^0 , that together compose the final metric given by Eq(10). The outputs of the network provide limited OOD detection capacity. Always when available, the intermediate features FR are a valuable resource for OOD detection. Moreover, when few reliable OOD data are available, calculating FR can further improve the detection performance as shown on the left side column of Table 1. Also, data from a

source other than in-distribution, e.g., adversarial samples, is enough for tuning hyperparameters and combining features, as observed on the right side column of Table 1. Figure 1 shows the detection performance for each hidden feature of a DenseNet as well as the scores histograms for a particular task for both Mahalanobis and IGEOOD scores.

(a) Block 1. (b) Block 2. (c) Block 3.

Figure 1: Histograms of the Mahalanobis and IGEOOD scores for the outputs of each hidden block of a DenseNet model on CIFAR-10 (in-distribution) and SVHN (out-of-distribution) datasets.

IGEOOD compared to other WHITE-BOX methods. Even though reference [20] shares the closest setup to ours, recent literature also proposes OOD detection in WHITE-BOX setting, achieving state-of-the-art in a few benchmarks. Notably, [15, 47] achieve great performance in a range of benchmarks. Thus, we report the results from the original references in Table 2. This setup considers that a few OOD samples are available for tuning. We refer the reader to the appendix (see Section D) where we provide a table of results with validation on adversarial data.

Table 2: TNR at TPR-95% (%) performance comparison in WHITE-BOX setting considering the original results from [20, 34, 15, 47]. Methods with an (*) are tuned only with in-distribution data.

	OOD dataset	CIFAR-10		CIFAR-100		SVHN	
		Mahalanobis [20] / Gram Matrix* [34]	DeConf-C* [15] / Res-Flow [47]	IGEOOD	IGEOOD+		
DenseNet	iSUN	95.3/99.0/ - / - /97.7/99.8	87.0/95.9/ - / - /93.8/99.7	99.9/99.4/ - / - /98.3/99.9	99.9/99.5/ - /100/97.1/99.9		
	LSUN	97.2/99.5/99.4/98.2/98.9/99.9	91.4/97.2/98.7/96.3/95.9/99.9	99.9/99.5/ - /100/97.1/99.9	99.9/99.1/ - /100/98.2/99.9		
	TinyImgNet	95.0/98.8/99.1/96.4/95.9/99.8	86.6/95.7/98.6/93.0/94.9/99.5	99.9/99.1/ - /100/98.2/99.9	99.9/99.1/ - /100/98.2/99.9		
	SVHN/C-10	90.8/96.1/98.8/94.9/98.9/99.9	82.5/89.3/95.9/84.9/93.9/99.6	96.8/80.4/ - /99.0/91.6/98.3	96.8/80.4/ - /99.0/91.6/98.3		
ResNet	iSUN	97.8/99.3/ - / - /97.7/99.9	89.9/94.8/ - / - /93.4/99.8	99.7/99.4/ - / - /99.8/100	99.7/99.4/ - / - /99.8/100		
	LSUN	98.8/99.6/ - /99.0/98.4/100	90.9/96.6/ - /96.2/94.3/100	99.9/99.6/ - /100/99.7/99.9	99.9/99.6/ - /100/99.7/99.9		
	TinyImgNet	97.1/98.7/ - /97.8/96.9/99.6	90.9/94.8/ - /94.6/90.9/99.6	99.9/99.3/ - /100/99.7/99.9	99.9/99.3/ - /100/99.7/99.9		
	SVHN/C-10	87.8/97.6/ - /96.5/98.9/99.8	91.9/80.8/ - /93.0/91.0/99.7	98.4/85.8/ - /99.4/97.9/99.7	98.4/85.8/ - /99.4/97.9/99.7		

4 Summary and Concluding Remarks

In this work, we introduced IGEOOD, an effective and flexible method for Out-Of-Distribution (OOD) detection that applies to any pre-trained neural network. The main feature of IGEOOD relies on the geodesic distance of the probabilistic manifold of the learned latent representations that induces an effective measure for OOD detection. The Fisher-Rao distance between pdf of the latent feature, corresponding to the test sample, and a reference pdf, corresponding to the conditional-class of pdfs, provides an effective confidence score. We consider diverse testing environments where prior knowledge of OOD data may or may not be available. Experimentally, we showed that IGEOOD can significantly and consistently improve the accuracy of OOD detection on several DNN models and across various OOD datasets, achieving new state-of-the-art performances on a few benchmarks.

Acknowledgments and Disclosure of Funding

This work has been supported by the project PSPC AIDA: 2019-PSPC-09 funded by BPI-France.

References

- [1] Dario Amodei, Chris Olah, Jacob Steinhardt, Paul F. Christiano, John Schulman, and Dan Mané. Concrete problems in AI safety. *CoRR*, abs/1606.06565, 2016.
- [2] Colin Atkinson and Ann F. S. Mitchell. Rao's distance measure. *Sankhyā: The Indian Journal of Statistics, Series A* (1961-2002), 43(3):345–365, 1981.
- [3] Abhijit Bendale and Terrance E. Boult. Towards open set deep networks. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1563–1572, 2016.
- [4] Julian Bitterwolf, Alexander Meinke, and Matthias Hein. Certifiably adversarially robust detection of out-of-distribution data. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 16085–16095. Curran Associates, Inc., 2020.
- [5] Michael M. Bronstein, Joan Bruna, Taco Cohen, and Petar Velicković. *Geometric deep learning: Grids, groups, graphs, geodesics, and gauges*, 2021.
- [6] Hyun-Jae Choi and Eric Jang. Generative ensembles for robust anomaly detection. *arXiv preprint arXiv:1810.01392*, 2018.
- [7] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, and A. Vedaldi. Describing textures in the wild. In *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- [8] T. E. de Campos, B. R. Babu, and M. Varma. Character recognition in natural images. *Proceedings of the International Conference on Computer Vision Theory and Applications, Lisbon, Portugal*, January 2009.
- [9] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009.
- [10] Ian Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*, 2015.
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- [12] Matthias Hein, Maksym Andriushchenko, and Julian Bitterwolf. Why relu networks yield high-confidence predictions far away from the training data and how to mitigate the problem. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 41–50, 2019.
- [13] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *International Conference on Learning Representations*, 2017.
- [14] Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. Deep anomaly detection with outlier exposure. In *International Conference on Learning Representations*, 2019.
- [15] 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 *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10948–10957, 2020.
- [16] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2261–2269, 2017.
- [17] Polina Kirichenko, Pavel Izmailov, and Andrew G Wilson. Why normalizing flows fail to detect out-of-distribution data. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 20578–20589. Curran Associates, Inc., 2020.
- [18] Alex Krizhevsky et al. Learning multiple layers of features from tiny images. 2009.
- [19] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. 2015.

- [20] Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, pages 7167–7177. Curran Associates, Inc., 2018.
- [21] Shiyu Liang, Yixuan Li, and R. Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. *International Conference on Learning Representations*, 2018.
- [22] Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection. *Advances in Neural Information Processing Systems*, 2020.
- [23] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *International Conference on Learning Representations*, 2018.
- [24] Prasanta Chandra Mahalanobis. On the generalized distance in statistics. *Proceedings of the National Institute of Sciences (Calcutta)*, 2:49–55, 1936.
- [25] Ahsan Mahmood, Junier Oliva, and Martin Andreas Styner. Multiscale score matching for out-of-distribution detection. *International Conference on Learning Representations*, 2021.
- [26] Sina Mohseni, Mandar Pitale, JBS Yadawa, and Zhangyang Wang. Self-supervised learning for generalizable out-of-distribution detection. *Proceedings of the AAAI Conference on Artificial Intelligence* 34(04):5216–5223, Apr. 2020.
- [27] Vaishnavh Nagarajan, Anders Andreassen, and Behnam Neyshabur. Understanding the failure modes of out-of-distribution generalization. *International Conference on Learning Representations*, 2021.
- [28] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y. Ng. Reading digits in natural images with unsupervised feature learning. *NIPS Workshop on Deep Learning and Unsupervised Feature Learning* 201, 2011.
- [29] Marco Pimentel, David Clifton, Lei Clifton, and L. Tarassenko. A review of novelty detection. *Signal Processing* 99:215–249, 06 2014.
- [30] Julianna Pinele, João E. Strapasson, and Sueli I. R. Costa. The χ^2 distance between multivariate normal distributions: Special cases, bounds and applications. *Entropy*, 22(4), 2020.
- [31] Joaquin Quionero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D. Lawrence. *Dataset Shift in Machine Learning*. The MIT Press, 2009.
- [32] Jie Ren, Stanislav Fort, Jeremiah Liu, Abhijit Guha Roy, Shreyas Padhy, and Balaji Lakshminarayanan. A simple χ^2 to mahalanobis distance for improving near-ood detection, 2021.
- [33] 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 H. Wallach, H. Larochelle, A. Beygelzimer, F. Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [34] Chandramouli Shama Sastry and Sageev Oore. Detecting out-of-distribution examples with Gram matrices. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning* volume 119 of *Proceedings of Machine Learning Research*, pages 8491–8501. PMLR, 13–18 Jul 2020.
- [35] Thomas Schlegl, Philipp Seeböck, Sebastian M. Waldstein, Ursula Margarethe Schmidt-Erfurth, and Georg Langs. Unsupervised anomaly detection with generative adversarial networks to guide marker discovery. In *IPMI*, 2017.
- [36] Gabi Shalev, Yossi Adi, and Joseph Keshet. Out-of-distribution detection using multiple semantic label representations. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- [37] João E. Strapasson, Julianna Pinele, and Sueli I. R. Costa. Clustering using the χ^2 distance. In *IEEE Sensor Array and Multichannel Signal Processing Workshop (SAMSP)*, pages 1–5, 2016.
- [38] Sachin Vernekar, Ashish Gaurav, Vahdat Abdelzad, Taylor Denouden, Rick Salay, and Krzysztof Czarnecki. Out-of-distribution detection in classifiers via generation. *Neural Information Processing Systems (NeurIPS 2019)*, *Safety and Robustness in Decision Making Workshop* <https://sites.google.com/view/neurips19-safe-robust-workshop>, <https://sites.google.com/view/neurips19-safe-robust-workshop>, 12/2019 2019.

- [39] 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. *ICCV* (8) pages 560–574, 2018.
- [40] Jim Winkens, Rudy Bunel, Abhijit Guha Roy, Robert Stanforth, Vivek Natarajan, Joseph R. Ledsam, Patricia MacWilliams, Pushmeet Kohli, Alan Karthikesalingam, Simon A. A. Kohl, taylan. cemgil, S. M. Ali Eslami, and Olaf Ronneberger. Contrastive training for improved out-of-distribution detection. *ArXiv*, abs/2007.05566, 2020.
- [41] J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. 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, 2010.
- [42] Zhisheng Xiao, Qing Yan, and Yali Amit. Likelihood regret: An out-of-distribution detection score for variational auto-encoder. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 20685–20696. Curran Associates, Inc., 2020.
- [43] Pingmei Xu, Krista A. Ehinger, Yinda Zhang, Adam Finkelstein, Sanjeev R. Kulkarni, and Jianxiong Xiao. Turkergaze: Crowdsourcing saliency with webcam based eye tracking. *ArXiv*, abs/1504.06755, 2015.
- [44] F. Yu, Y. Zhang, Shuran Song, Ari Seff, and J. Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *ArXiv*, abs/1506.03365, 2015.
- [45] Yufeng Zhang, Wanwei Liu, Zhenbang Chen, Ji Wang, Zhiming Liu, Kenli Li, and Hongmei Wei. Towards out-of-distribution detection with divergence guarantee in deep generative models. 2020.
- [46] 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* 2017.
- [47] Ev Zisselman and Aviv Tamar. Deep residual flow for out of distribution detection. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

A Review of Fisher-Rao Distance (FRD)

In this section, we review some results from reference [2]. We intend to clarify some basic concepts surrounding the Fisher-Rao distance while motivating the use of this measure in the context of OOD detection.

In few words, the Fisher-Rao's distance is given by the geodesic distance, i.e., the shortest path between points in a Riemannian space induced by a parametric family. Consider the family of probability distributions over the class of discrete concepts or labels $\mathcal{Y} = \{1, \dots, C\}$, denoted by \mathcal{C} , $q(y|x) : x \in \mathcal{X} \rightarrow \mathcal{R}^C$.

We are interested in measuring the distance between probability distributions with respect to the testing input x and a population of inputs drawn accordingly to the in-distribution data set. To this end, we first need to characterize the Fisher-Rao distance for two inputs or for two probability distributions \mathcal{C} .

Assume that the following regularity conditions hold [2]:

- (i) $r_x q(y|x)$ exists for all x, y and \mathcal{Y} ;
- (ii) $r_x r_y q(y|x) = 0$ for all x and \mathcal{Y} ;
- (iii) $G(x) = E_{q(y|x)} [r_x \log q(Y|x) r_x^T \log q(Y|x)]$ is positive definite for any x and \mathcal{Y} .

Notice that if (i) holds, (ii) also holds immediately for discrete distributions over finite spaces (assuming that $r_x r_y$ and r_x are interchangeable operations) as in our case. When (i)-(iii) are met, the variance of the differential form $r_x \log q(Y|x) dx$ can be interpreted as the square of a differential arc length in the space \mathcal{C} , which yields

$$ds^2 = dx^T G(x) dx \quad (11)$$

Thus, G , which is the Fisher Information Matrix (FIM), can be adopted as a metric tensor. We now consider a curve $\gamma : [0, 1] \rightarrow \mathcal{X}$ connecting a pair of arbitrary points x, x^0 in the input space \mathcal{X} , i.e., $\gamma(0) = x$ and $\gamma(1) = x^0$. Notice that any curve induces a curve $q(\gamma(t))$ for $t \in [0, 1]$ in the space \mathcal{C} . The Fisher-Rao distance between the distributions $q = q(\gamma(x))$ and $q^0 = q(\gamma(x^0))$ will be denoted as $d_{R,C}(q; q^0)$ and is formally defined by the expression:

$$d_{R,C}(q; q^0) = \inf_{\gamma} \int_0^1 \sqrt{\frac{d}{dt} \gamma^T G(\gamma(t)) \frac{d}{dt} \gamma} dt \quad (12)$$

where the inimum is taken over all piecewise smooth curves. This means that the FRD is the length of the geodesic between points x and x^0 using the FIM as the metric tensor. In general, the minimization of the functional in Eq.(12) is a problem that can be solved using the well-known Euler-Lagrange differential equation.

A.1 Derivation of Fisher-Rao distance for the class of Softmax probability distributions

The direct computation of the FIM of the family with $q(y|x)$ in the form of the softmax probability distribution function Eq.(1) can be shown to be singular, i.e., $\text{rank}(G(x)) < C - 1$, where $C - 1$ is the number of degrees of freedom of the manifold \mathcal{C} . To overcome this issue, we introduce the probability simplex defined by

$$P = \{q : \mathcal{Y} \rightarrow [0, 1]^C : \sum_{y \in \mathcal{Y}} q(y) = 1\} \quad (13)$$

Next, we consider the following parametrization for any distribution $q \in P$:

$$q(y|z) = \frac{z_y^2}{4}; \quad y \in \{1, \dots, C\} \quad (14)$$

From this expression, we consider the statistical manifold $\mathcal{D} = \{q(z) : \|z\|^2 = 4; z_y \geq 0; y \in \mathcal{Y}\}$. Note that the parameter vector z belongs to the positive portion of a sphere of radius 2 and centered at the origin in \mathcal{R}^C . The computation of the FIM for \mathcal{D} yields:

$$\begin{aligned} G(z) &= E_{q(y|z)} [r_z \log q(y|z) r_z^T \log q(y|z)] \\ &= \sum_{y \in \mathcal{Y}} \frac{z_y^2}{4} \frac{2}{z_y} e_y e_y^T \\ &= \sum_{y \in \mathcal{Y}} e_y e_y^T \\ &= I; \end{aligned} \quad (15)$$

where e_y, g are the canonical basis vectors in \mathbb{R}^n and I is the identity matrix. From expression (15) we can conclude that the Fisher-Rao metric in this parametric space is equal to the Euclidean metric. Also, since the parameter vector lies on a sphere, the FRD between the distributions $q(z)$ and $q^0 = q(z^0)$ can be written as the radius of the sphere times the angle between the vectors z and z^0 . Which leads to expression:

$$d_{R; D} (q; q^0) = 2 \arccos \frac{z^T z^0}{\|z\| \|z^0\|} = 2 \arccos \frac{\sum_{y=1}^n p_y \frac{q(y|z)q(y|z^0)}{\|q(y|z)\| \|q(y|z^0)\|}}{\|q(y|z)\| \|q(y|z^0)\|} \quad (16)$$

Finally, we can compute the FRD for softmax distribution $q(x)$ as

$$d_{FR; Logits} (q; q^0) = 2 \arccos \frac{\sum_{y=1}^n p_y \frac{q(y|x)q(y|x^0)}{\|q(y|x)\| \|q(y|x^0)\|}}{\|q(y|x)\| \|q(y|x^0)\|}; \quad (17)$$

obtaining the same form of expression (16). Notice that $d_{FR; Logits} (q; q^0) = 0$ for all $x; x^0 \in \mathbb{R}^C$, being zero when $q(y|x) = q(y|x^0)$ and maximum when the vectors $q(1|x); \dots; q(C|x)$ and $q(1|x^0); \dots; q(C|x^0)$ are orthogonal.

A.2 Derivation of Fisher-Rao distance for multivariate Gaussian distributions

Consider a broader statistical manifold \mathcal{M} , $f_p = p(x; \theta) = \frac{1}{\sqrt{|P_k(\theta)|}} \exp(-\frac{1}{2}(x - \mu)^T P_k^{-1}(\theta)(x - \mu))$ of multivariate differential probability density functions. The Fisher information matrix $g_{ij}(\theta) = \int \frac{\partial}{\partial \theta_i} \log p(x; \theta) \frac{\partial}{\partial \theta_j} \log p(x; \theta) p(x; \theta) dx$ in this parametric space is provided by:

$$g_{ij}(\theta) = E \left[\frac{\partial}{\partial \theta_i} \log p(x; \theta) \frac{\partial}{\partial \theta_j} \log p(x; \theta) \right] = \int \frac{\partial}{\partial \theta_i} \log p(x; \theta) \frac{\partial}{\partial \theta_j} \log p(x; \theta) p(x; \theta) dx \quad (18)$$

Next, consider a multivariate Gaussian distribution:

$$p(x; \mu, P_k) = \frac{1}{\sqrt{|P_k|}} \exp\left(-\frac{1}{2}(x - \mu)^T P_k^{-1}(x - \mu)\right); \quad (19)$$

where $x \in \mathbb{R}^k$ is the variable vector, $\mu \in \mathbb{R}^k$ is the mean vector, $P_k \in \mathbb{R}^{k \times k}$ is the covariance matrix, and $P_k \in \mathbb{R}^{k \times k}$ is the space of positive definite symmetric matrices. We can define the statistical manifold composed by these distributions $\mathcal{M} = \{p; \mu \in \mathbb{R}^k, P_k \in \mathbb{R}^{k \times k}\}$. By substituting Eq(19) in expression Eq. (18), we can derive the Fisher information matrix for this parametrization, obtaining:

$$g_{ij}(\mu, P_k) = \frac{\partial \mu^T}{\partial \theta_i} \frac{\partial \mu}{\partial \theta_j} + \frac{1}{2} \text{tr} \left[\frac{\partial P_k^{-1}}{\partial \theta_i} \frac{\partial P_k^{-1}}{\partial \theta_j} \right]; \quad (20)$$

which induces the following square differential arc length ds^2 :

$$ds^2 = d\mu^T d\mu + \frac{1}{2} \text{tr} \left[\frac{\partial P_k^{-1}}{\partial \theta_i} \frac{\partial P_k^{-1}}{\partial \theta_j} \right] d\theta_i d\theta_j; \quad (21)$$

Here, $d = (d_1; \dots; d_n) \in \mathbb{R}^k$ and $d = [d_{ij}] \in \mathbb{R}^{k \times k}$. We observe that this metric is invariant to affine transformations [30], i.e., for any $(c; Q) \in \mathbb{R}^k \times GL_k(\mathbb{R})$, with $GL_k(\mathbb{R})$ the space of non-singular order k matrices, the map $(\mu; P_k) \mapsto (c + Q\mu; QP_kQ^T)$ is an isometry in \mathcal{M} . Thus, the Fisher-Rao distance between two multivariate normal distributions with parameters $\theta_1 = (\mu_1; P_{k1})$ and $\theta_2 = (\mu_2; P_{k2})$ in \mathcal{M} satisfies:

$$d_{R; M}(\theta_1; \theta_2) = d_{R; M}(\mu_1 + c; Q^{-1}P_{k1}Q^{-T}; \mu_2 + c; Q^{-1}P_{k2}Q^{-T}); \quad (22)$$

Unfortunately, a closed-form solution for the Fisher-Rao distance remains unknown. This is still an open problem for an arbitrary covariance matrix and mean vector. Fortunately, the FRD is known for the univariate case and hence, for the submanifold where P_k is diagonal. Notice that in this case Eq. (21) admits an additive form.

From [30], we obtain the analytical expression of the Fisher-Rao in the 2-dimensional submanifold of univariate Gaussian probability distributions $\mathcal{M}_2 = \{p; \mu = (\mu_1; \mu_2) \in \mathbb{R}^2, \sigma_i^2 > 0; i = 1, 2\}$:

$$d_{FR; \text{Gauss}}(\mu_1; \sigma_1^2; \mu_2; \sigma_2^2) = \sqrt{2} \log \frac{\sqrt{\sigma_1^2 \sigma_2^2} \left(\frac{\mu_1^2}{\sigma_1^2} + \frac{\mu_2^2}{\sigma_2^2} + \frac{\mu_1^2}{\sigma_1^2} + \frac{\mu_2^2}{\sigma_2^2} \right)}{\sqrt{\sigma_1^2 \sigma_2^2} \left(\frac{\mu_1^2}{\sigma_1^2} + \frac{\mu_2^2}{\sigma_2^2} + \frac{\mu_1^2}{\sigma_1^2} + \frac{\mu_2^2}{\sigma_2^2} \right)}; \quad (23)$$

where $\|j\|$ is the Euclidean norm in \mathbb{R}^2 and σ_i denotes the standard deviation. Consequently, the FRD for Gaussian distributions with diagonal covariance matrix $P_k = \text{diag}(\sigma_1^2; \sigma_2^2; \dots; \sigma_k^2)$ in the $2k$ -dimensional statistical submanifold $\mathcal{M}_D = \{p; \mu = (\mu_1; \mu_2); \sigma_i^2 > 0; i = 1; \dots; k\}$ is

$$d_{FR; \text{Gauss}}(\mu_1; \mu_2) = \sqrt{\sum_{i=1}^k \frac{\sigma_i^2}{2} \left(\frac{\mu_{1i} - \mu_{2i}}{\sigma_i} \right)^2}; \quad (24)$$

A.3 Fisher-Rao vs. Mahalanobis distance

There is an intricate relationship between the FRD for multivariate Gaussian distributions and the Mahalanobis distance. We borrow the result from [30], which states that in the k -dimensional submanifold of M where μ is constant, i.e. $M = \{p : \mu = (\mu_1, \dots, \mu_k)\}$, the Fisher-Rao distance $d_{R;M}$ between two distributions is given by the Mahalanobis distance [24]:

$$d_{R;M}(N(\mu_1; \Sigma); N(\mu_2; \Sigma)) = \sqrt{(\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 - \mu_2)} \quad (25)$$

The Mahalanobis distance is also used for OOD detection and its performance is compared to the FRD through several experiments in Section 3.

B IGEOD Algorithms

In this section, we provide pseudo-code for calculating the IGEOD score from the logits (Algorithm 1) and from the latent features (Algorithm 2). The BLACK-BOX IGEOD score is obtained with Algorithm 1 by setting $\alpha = 0$, while the GREY-BOX IGEOD score is obtained with $\alpha > 0$. We calculated the centroid of the logits for the in-distribution training set by optimizing the objective function given by (9) through a gradient descent algorithm for each DNN. We used a constant learning rate of 0.01 and a batch size of 28 for 100 epochs. Finally, the WHITE-BOX IGEOD score is obtained by combining the outputs of Algorithms 1 and 2 through fitting the multiplicative weights through a logistic function classifier on a labeled mixture dataset composed from in- and out-of-distribution data according to a validation dataset, which leads to expression Eq. (10).

Algorithm 1: Evaluating IGEOD score based on the logits.

Input : Test sample x , temperature T and noise magnitude σ parameters, and training set

$$D_N = f(x_i; y_i)_{i=1}^N$$

Output : FR_0 : IGEOD score in the logits level.

// Offline computation

Calculate the logits centroids from the training data:

$$y_c = \arg \min_{y \in \mathcal{R}^c} \frac{1}{N_y} \sum_{i: y_i=y} 2 \arccos \frac{\sum_{i: y_i=y} q(y_i) f(x_i) q(y_i)}{q(y)^2}$$

// Online computation

Add small perturbation α :

$$x = x + \alpha \cdot \text{sign} \left(\sum_{i: y_i=y} 2 \arccos \frac{\sum_{i: y_i=y} q(y_i) f(x_i) q(y_i)}{q(y)^2} - y_c \right)$$

return $FR_0(x) = \sum_{i: y_i=y} 2 \arccos \frac{\sum_{i: y_i=y} q(y_i) f(x_i) q(y_i)}{q(y)^2}$

Algorithm 2: Evaluating feature-wise IGEOD score.

Input : Test sample x and training set $D_N = f(x_i; y_i)_{i=1}^N$.

Output : FR_j : feature-wise IGEOD scores.

for each feature $j \in \{1, \dots, L\}$ do

// Offline computation

Calculate the means: $\mu_j^{(y)} = \frac{1}{N_y} \sum_{i: y_i=y} f_j^{(y)}(x_i)$

Calculate the diagonal standard deviation matrix:

$$\Sigma_{jj}^{(y)} = \frac{1}{N_y} \sum_{i: y_i=y} f_j^{(y)}(x_i)^2 - \mu_j^{(y)2}$$

// Online computation

Compute the OOD score for

$$FR_j(x) = \min_y \sum_{k=1}^K \frac{FR_j^{(y)}(x)}{\Sigma_{jj}^{(y)}}; \quad f_j^{(y)}(x); \quad \Sigma_{jj}^{(y)2}$$

end

return $FR_1(x); \dots; FR_L(x)$

Algorithm 3: Evaluating feature-wise EOOD+ score.

Input : Test samples \mathcal{X} , training set $\mathcal{D}_N = \{f(x_i; y_i)\}_{i=1}^N$ and M OOD samples

$$\mathcal{O}_M = \{f(x_i^0; y_i^0)\}_{i=1}^M.$$

Output: FR_j and FR_j^0 : feature-wise EOOD+ scores.

for each feature $j \in \{1, \dots, L\}$ do

// Offline computation

Calculate class conditional means $\mu_j^{(c)}$ $\frac{1}{N_y} \sum_{i: y_i = y} f^{(c)}(x_i)$

Calculate OOD samples mean $\mu_j^{(0)}$ $\frac{1}{M} \sum_{i=1}^M f^{(c)}(x_i^0)$

Calculate the diagonal standard deviation matrix from training data:

$$\Sigma_{jj}^{(c)} = \frac{1}{N} \sum_{i: y_i = y} (f_j^{(c)}(x_i) - \mu_j^{(c)})^2$$

Calculate the diagonal standard deviation matrix from OOD data:

$$\Sigma_{jj}^{(0)} = \frac{1}{M} \sum_{i=1}^M (f_j^{(c)}(x_i^0) - \mu_j^{(0)})^2$$

// Online computation

Compute the OOD scores for

$$FR_j(x) = \min_y \sqrt{\frac{\sum_{j=1}^k FR_j^{(c)}(\mu_j^{(c)}; \mu_j^{(c)}; f_j^{(c)}(x); \Sigma_{jj}^{(c)})^2}{r}}$$

$$FR_j^0(x) = \min_y \sqrt{\frac{\sum_{j=1}^k FR_j^{(0)}(\mu_j^{(0)}; \mu_j^{(0)}; f_j^{(c)}(x); \Sigma_{jj}^{(0)})^2}{r}}$$

end

return $FR_1(x); FR_1^0(x); \dots; FR_L(x); FR_L^0(x)$

Note that the calculation of the training logits centroids $\mu_j^{(c)}$ as well as the latent representations mean vectors $\mu_j^{(0)}$ and standard covariance matrices $\Sigma_{jj}^{(c)}$ is performed beforehand, prior to inference. In this way, we retrieve the objects from memory at inference time. Also, we denote the cardinality of feature j as FR_j as the Fisher-Rao distance between univariate Gaussian distribution given by expression Eq. (6).

B.1 Logits centroids estimation details

In order to obtain the logits centroids given the Fisher-Rao distance in the space of softmax probability distributions, we designed a simple optimization problem. This problem aims to minimize the average distance between the class conditional training samples and the centroids as given in Eq. (9). We initialized the centroids, where C is the number of classes of a given model, with the identity matrix of size C . Note that the initial centroid for class c is given by the matrix's line number c . We minimized the expression in Eq. (9) with a gradient descent optimizer for 100 epochs with a fixed learning rate equal to 0.01 for every DNN model and in-distribution dataset.

B.2 Feature importance regression details

For both Mahalanobis and EOOD methods, we fitted a logistic regression model with cross-validation using 1,000 OOD and 1,000 in-distribution data samples. Each regression parameter multiplies the layer scores outputs with the objective function of maximizing the TNR at TPR-95%. We set the maximum number of iterations to be 100.

B.3 Covariance matrix estimation details

We model the latent output probability distributions as Gaussian distributions with diagonal covariance matrix calculated with expression Eq. (5). We chose this model motivated by a closed form for the FRD and by observing that the standard covariance matrix for the latent features is often ill-conditioned. The condition number of a matrix correlates to its numerical stability, i.e., a small rounding error in its estimation may cause a large difference in its values. So, a matrix with a low condition number is said to be well-conditioned, while a matrix with a high condition number is said to be ill-conditioned. We calculate the condition number of

the covariance matrices with the formula $\hat{\Sigma} = \frac{1}{n} \sum_{k=1}^n k_k k_k^T$, where k_k is the infinity norm. For each of the four dense blocks outputs of a DenseNet trained on CIFAR-10, we obtained the condition numbers $\kappa = [2.8e10, 3.5e6, 3.1e5, 3.5e2]g$. While for the diagonal covariance matrix, we obtained smaller values of condition numbers: $\kappa_D = [1.0e3, 3.0e1, 1.4e1, 7.6e2]g$. We associate the high value for the last feature due to the last feature being high dimensional and coarse, i.e., most of the values in the diagonal are close to zero.

C Detailed Experimental Setup

C.1 DNN models and training details

We describe the DNN models used in the experiments:

- **DenseNet**. Densely Connected Convolutional Networks^[6], or DenseNet for short, are compositions of dense blocks, which are composed of multiple layers directly connected to every other layer in a feed-forward fashion. In this work, we use the DenseNet-BC-100 architecture. The BC stands for a model with 1x1 convolutional bottleneck (B) layers and channel number compression (C) of 0.5. The models have $depth = 100$ and $growth\ rate = 12$. We consider the outputs of each dense block after the transition layer (3 in total) and the first convolutional layer output as the latent features. After an averaging pooling, the latent features have dimensions $F_2 = [24; 108; 150; 342]g$.
- **ResNet**. Residual Networks^[11], or ResNet, are deep neural networks composed of residual blocks. Each residual block is composed of layers connected in a feed-forward manner plus a skip connection. We use the ResNet with 34 layers pre-trained on CIFAR-10, CIFAR-100, and SVHN datasets. We take the output of every residual block (4 in total) and the first convolutional layer for calculating the score on the WHITE-BOX setting. After an averaging pooling, the latent features have dimensions $F_2 = [64; 64; 128; 256; 512]g$.

We train each model by minimizing the cross-entropy loss using SGD with Nesterov momentum equal to 0.9, weight decay equal to 0.0001, and a multi-step learning rate schedule starting at 0.1 for 300 epochs. The pre-trained models is available³ at We report their test set accuracy in Table 3 with the softmax function and by replacing it with the Fisher-Rao distance between the training class-conditional centroids and the test sample outputs. Also, it is worth noting that one high-end GPU is sufficient for running every experiment presented in this work.

Table 3: Test set accuracy in percentage for ResNet and DenseNet architectures pre-trained on CIFAR-10, CIFAR-100 and SVHN.

In-Dataset	ResNet-34		DenseNet-BC-100	
	Softmax	Fisher-Rao	Softmax	Fisher-Rao
CIFAR-10	93.52	93.53	95.20	95.20
CIFAR-100	77.11	77.09	77.62	77.63
SVHN	96.61	96.61	95.16	95.16

C.2 Evaluation metrics

We introduce below standard binary classification performance metrics used to evaluate the OOD discriminators.

- **True Negative Rate at 95% True Positive Rate (TNR at TPR-95% (%))**. This metric measures the true negative rate (TNR) at a specific true positive rate (TPR). The operating point is chosen such that the TPR of the in-distribution test set is fixed to some value, 95% in this case. Mathematically, let TP, TN, FP, and FN denote true positive, true negative, false positive and false negative, respectively. We measure $TNR = \frac{TN}{(FP + TN)}$, when $TPR = \frac{TP}{(TP + FN)}$ is 95%.
- **Area Under the Receiver Operating Characteristic curve (AUROC (%))**. The ROC curve is constructed by plotting the true positive rate (TPR) against the false positive rate (FP+ TN) at various threshold values. The area under this curve tells how much the OOD discriminator can distinguish in-distribution and OOD data in a threshold-independent manner.
- **Area Under the Precision-Recall curve (AUPR (%))**. The PR curve plots the precision $\frac{TP}{(TP + FN)}$ against the recall $\frac{TP}{(TP + FN)}$ by varying a threshold. For the experiments, in-distribution data are specified as positives while OOD data as negative.

³<http://github.com/igeood/lgeood>

Note that the TNR at TPR-95% is especially important because we want to identify OOD data and preserve a sufficiently good performance on identifying in-distribution data, which is not the case for the other metrics.

C.3 Datasets

In our experiments, we use natural image examples from the following image classification and synthetic datasets. We normalize the test samples with the in-distribution dataset statistics.

- CIFAR-10. The CIFAR-10 [18] dataset is composed of 32 × 32 natural images of 10 different classes, e.g., airplane, ship, bird, etc. The training set is composed of 50,000 images, and the test set is composed of 10,000 images. The classes are approximately equally distributed (5,000 examples each label). The CIFAR-10 dataset is under the MIT license.
- CIFAR-100. The CIFAR-100 [18] dataset contains similar natural images to the CIFAR-10 dataset, but with 90 additional categories. Its set repartition is also 50,000 for training and 10,000 for the test set. We expect around 500 samples for each class of the training set. It is also under the MIT license.
- SVHN. The SVHN [28] dataset collects street house numbers for digit classification. It contains 73,257 training and 26,032 test RGB images of size 32 × 32 of printed digits (from 0 to 9). We take only the first 10,000 examples of the test set for evaluating the methods to have a balanced dataset of in-distribution and out-of-distribution data. This dataset is subject to a non-commercial license.
- Tiny-ImageNet. The Tiny-ImageNet [9] dataset is a subset of the large-scale natural image dataset ImageNet [9]. It contains 200 different classes and 10,000 test examples. We downsize the images from their original resolution to images of dimension 32 × 32.
- LSUN. The LSUN [44] dataset, which has equally 10,000 test examples, is used for the large-scale scene classification of different scene categories (e.g., bedroom, bridge, kitchen, etc.). Similarly, we resize the images following the same procedure for the Tiny-ImageNet dataset. LSUN is under the Apache 2.0 license.
- iSUN. The iSUN [43] dataset consists of selected natural scene images from the SUN dataset. The test set has 8925 images, which we downsample to 32 × 32. We use this dataset as a source of OOD for validation purposes as an independent dataset from the test OOD data.
- Textures. The Describable Textures Dataset (DTD) [5] is a collection of textural pattern images observed in nature. It contains 47 categories totaling 5640 images of various sizes, which are resized and center cropped to fit into the input size 32 × 32.
- Chars74K. The Chars74K dataset [6] contains 74,000 samples of 62 classes of characters found in natural images, handwritten text, and synthesized from computer fonts. We used as OOD data only the EnglishImg dataset split, which contains 7705 characters from natural scenes. We resized and center-cropped the images.
- Places365. The Places365 dataset [46] contains images of 365 natural scenes categories. We used the small images validation split as OOD data in our experiments. It contains 36,500 RGB images which were downsampled from 256 × 256 to 32 × 32.
- Gaussian. For the Gaussian dataset, we generated 10,000 synthetic RGB images from 2D Gaussian noise, where each RGB pixel is sampled from an i.i.d Gaussian distribution with mean 0.5 and variance 1.0. The pixel values are clipped to [0, 1] interval. This synthetic data was introduced in previous work as an easy benchmark [13].

C.4 Adversarial data generation

We generate adversarial samples from the in-distribution dataset using the fast gradient sign method (FGSM). This method works by exploiting the gradients of the neural network to create a non-targeted adversarial attack. For an input image x_i , the method computes the sign of the gradients of the loss function with respect to the input image to create a new image x_i^{adv} that maximizes the loss as given by expression (26). This fabricated image is called an adversarial image, which we use for tuning the hyperparameters of the OOD detection methods in the MTE-BOX case. Mathematically,

$$x_i^{adv} = x_i + \epsilon^{adv} \text{sign}(\nabla_{x_i} J(\cdot; x_i; y_i)); \quad (26)$$

where $\epsilon^{adv} > 0$ is the additive noise magnitude parameter. Table 4 shows the resulting mean perturbation and classification accuracy on adversarial samples.

C.5 Mahalanobis distance-based confidence score

The Mahalanobis-based method [20] treats the DNN training data features as class-conditional Gaussian distributions. These use the outputs of every DNN latent block to leverage useful information for discrimination.

Table 4: The L_1 mean perturbation used to generate adversarial data with FGSM algorithm and classification accuracy on adversarial samples for the DNN models and in-distribution datasets.

	CIFAR-10		CIFAR-100		SVHN	
	L_1	Acc.	L_1	Acc.	L_1	Acc.
DenseNet-100	0.21	19.5%	0.20	4.45%	0.32	54.7%
ResNet-34	0.21	23.7%	0.20	12.49%	0.25	50.0%

For a test sample x , the confidence score from the l -th feature is calculated based on the Mahalanobis distance between $f^{(l)}(x)$ and the closest class-conditional distribution:

$$M_l(x) = \max_y \left(f^{(l)}(x) - b_y^{(l)} \right)^T \Sigma_y^{-1} \left(f^{(l)}(x) - b_y^{(l)} \right); \quad (27)$$

where $f^{(l)}(x)$ is the l -th latent feature output, $a^{(l)}$ and $b_y^{(l)}$ are, respectively, the empirical class mean and covariance matrix estimates. The covariance matrix is often not full rank, so the pseudo-inverse is calculated instead of the inverse. In addition, input pre-processing and feature ensemble are also used to boost performance. A logistic regression model learns the multiplicative weights for each layer score, which predicts for in-distribution and OOD for OOD examples from a mixture validation dataset. Finally, the Mahalanobis-based discriminator is given by thresholding expression $M_l(x)$.

D Additional Out-Of-Distribution detection results

Table 5 shows comparisons to the current literature in the setup where no OOD data is available for tuning. We show in Tables 6 and 7 additional results referring to the right-hand column and left-hand column of Table 1, respectively.

Table 5: TNR at TPR-95% (%) performance in a WHITE-BOX setting considering the original results from [20, 34, 15, 47] without access to OOD samples for hyperparameter tuning.

	OOD dataset	CIFAR-10		CIFAR-100		SVHN	
		Mahalanobis	Gram Matrix	DeConf-C	Res-Flow	IGOOD+	IGOOD+
DenseNet	iSUN	94.3/99.0	99.4/94.5	84.8/95.9	98.4/93.8	99.9/99.4	99.2/98.2
	LSUN	97.2/99.5	99.4/98.1	91.4/97.9	97.7/95.8	100/99.5	100/97.3
	TinyImgNet	94.9/98.8	99.1/96.1	87.2/95.9	98.6/91.5	99.9/99.1	99.9/98.1
	SVHN/C-10	89.9/96.8	88.6/94.3	62.2/89.9	95.9/48.9	90.0/80.4	90.0/89.5
	average	94.1/98.9	99.2/93.4	81.4/94.5	97.9/78.7	97.4/94.6	99.6/95.8
ResNet	iSUN	96.8/99.3	88.8/95.3	87.9/84.8	75.3/89.4	100/99.4	99.8/99.9
	LSUN	98.1/99.6	90.9/99.1	56.9/67.6	87.8/70.4	99.9/99.6	100/99.8
	TinyImgNet	95.5/98.7	81.4/98.0	70.3/84.8	76.5/77.5	99.2/99.3	99.9/99.6
	SVHN/C-10	75.8/97.6	89.5/91.8	41.9/80.8	85.1/74.1	94.1/85.8	96.6/96.7
	average	91.5/88.8	87.6/96.0	64.9/71.0	74.0/84.8	98.3/96.0	99.0/99.2

E Histograms

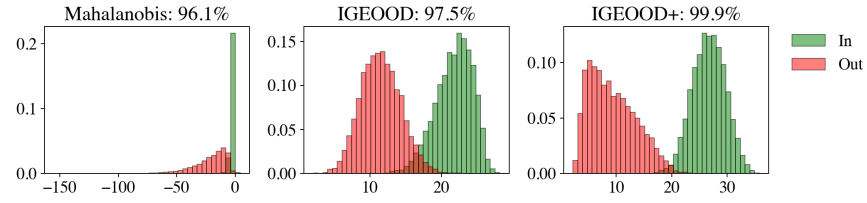
Figures 3 and 2 display histograms for the OOD detection scores for IGOD and [20] in the WHITE-BOX with adversarial validation and WHITE-BOX with OOD data validation, respectively.

Table 6: WHITE-BOX extended results. Validation on OOD data.

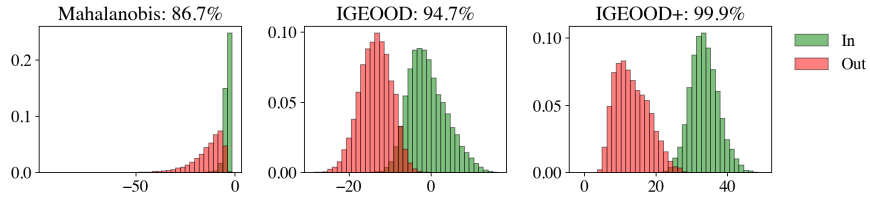
In-dist. (model)	OOD dataset	TNR at TPR-95%	AUROC	AUPR
		Mahalanobis [20] / \mathbb{E}_{OOD+}		
CIFAR-10 (DenseNet)	Chars	91.399.4	97.599.9	97.799.9
	CIFAR-100	21.456.6	67.390.7	64.490.8
	TinyImgNet	96.999.8	99.399.9	99.399.9
	LSUN	98.299.9	99.5100	99.5100
	Places365	18.80.2	72.795.7	72.895.4
	SVHN	90.199.9	97.3100	97.3100
	Textures	84.97.4	95.699.5	94.799.5
	Gaussian	100100	100100	100100
	iSUN	97.399.8	99.4100	99.4100
	average	77.531/92.6 ¹⁴	92.1 ¹² /98.4 ^{3.0}	91.7 ¹³ /98.4 ^{3.0}
CIFAR-100 (DenseNet)	Chars	62.997.5	94.099.4	95.899.4
	CIFAR-10	9.122.7	60.880.7	60.183.0
	TinyImgNet	87.199.5	97.499.9	97.499.9
	LSUN	91.199.9	97.8100	98.1100
	Places365	5.98.2	54.890.0	54.789.2
	SVHN	79.099.6	96.899.9	94.199.9
	Textures	70.390.2	91.498.1	94.398.2
	Gaussian	100100	100100	100100
	iSUN	86.499.7	96.899.9	97.799.9
	average	67.728/90.2 ²¹	87.8 ¹³ /97.7 ^{5.0}	88.0 ¹² /97.8 ^{5.0}
SVHN (DenseNet)	Chars	78.792.2	96.198.4	98.998.5
	CIFAR-10	91.698.3	98.099.6	99.499.6
	CIFAR-100	92.995.3	98.299.1	99.499.2
	TinyImgNet	99.999.9	99.899.9	99.999.9
	LSUN	99.999.9	99.8100	99.7100
	Places365	94.98.3	98.399.6	98.499.7
	Textures	98.298.5	99.499.6	99.999.6
	Gaussian	100100	100100	100100
	iSUN	99.999.9	99.899.9	99.999.9
	average	95.180/98.0 ^{2.0}	98.8 ^{1.0} /99.6 ^{0.1}	99.5 ^{1.0} /99.6 ^{0.1}
CIFAR-10 (ResNet)	Chars	93.699.3	98.699.8	99.199.8
	CIFAR-100	44.951.3	87.490.9	87.891.7
	TinyImgNet	96.899.6	99.499.9	99.499.9
	LSUN	98.399.9	99.6100	99.6100
	Places365	45.877.6	88.195.6	88.195.5
	SVHN	96.199.8	99.099.9	98.199.9
	Textures	84.397.0	97.399.4	98.699.4
	Gaussian	100100	100100	100100
	iSUN	97.299.9	99.4100	99.5100
	average	84.123/91.6 ¹⁶	96.5 ^{4.0} /98.4 ^{3.0}	96.7 ^{4.0} /98.5 ^{3.0}
CIFAR-100 (ResNet)	Chars	63.897.8	94.099.5	96.099.5
	CIFAR-10	18.080.8	76.685.3	76.487.8
	TinyImgNet	90.199.6	97.999.9	98.099.9
	LSUN	92.4100	98.3100	98.5100
	Places365	23.59.1	76.891.2	76.091.4
	SVHN	88.499.7	97.799.9	95.299.9
	Textures	71.690.7	93.998.2	96.698.1
	Gaussian	100100	100100	100100
	iSUN	89.499.8	97.799.9	98.099.9
	average	70.830/86.4 ²³	92.5 ¹⁰ /97.1 ^{5.0}	92.7 ¹⁰ /97.4 ^{4.0}
SVHN (ResNet)	Chars	84.992.4	97.098.4	99.098.5
	CIFAR-10	98.099.7	99.299.9	99.799.9
	CIFAR-100	98.399.1	99.399.7	99.899.8
	TinyImgNet	99.999.9	99.9100	100100
	LSUN	99.999.9	99.9100	100100
	Places365	98.89.6	99.399.9	99.899.9
	Textures	99.099.9	99.799.9	99.999.9
	Gaussian	100100	100100	100100
	iSUN	100100	99.9100	100100
	average	97.660/98.9 ^{2.0}	99.4 ^{1.0} /99.7 ^{0.1}	99.8 ^{1.0} /99.8 ^{0.1}
Avg. and std. of avg. values	82.11/92.9 ^{4.0}	94.5 ^{4.0} /98.5 ^{1.0}	94.7 ^{4.0} /98.6 ^{1.0}	

Table 7: WHITE-Box extended results. Validation on adversarial (FGSM) data.

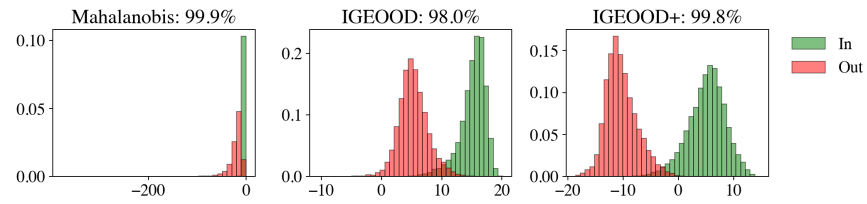
In-dist. (model)	OOD dataset	TNR at TPR-95%	AUROC	AUPR
		Mahalanobis [20] / ϵ_{OOD}		
CIFAR-10 (DenseNet)	Chars	88.587.3	97.797.7	98.398.3
	CIFAR-100	21.526.4	68.077.7	66.375.5
	TinyImageNet	93.993.4	98.698.7	98.698.7
	LSUN	96.396.4	99.199.2	99.199.2
	Places365	17.88.2	70.077.9	40.176.3
	SVHN	87.094.3	97.298.7	93.797.3
	Textures	83.66.0	95.897.2	97.398.2
	Gaussian	100100	100100	100100
	iSUN	94.394.5	98.898.9	98.999.0
	average	75.930777.9 ²⁹	91.7 ¹² 94.0 ^{9.0}	88.0 ²⁰ 93.6 ¹⁰
CIFAR-100 (DenseNet)	CIFAR-10	1.15.7	43.562.6	46.762.6
	Chars	53.979.6	92.292.0	94.593.9
	TinyImageNet	86.494.3	97.498.8	97.598.9
	LSUN	88.695.1	97.698.9	97.998.9
	Places365	5.53.0	56.671.0	57.571.0
	SVHN	56.190.1	91.898.0	85.496.2
	Textures	67.596.7	91.297.4	94.498.4
	Gaussian	100100	100100	100100
	iSUN	84.893.8	97.298.7	97.698.8
	average	60.43470.9 ³⁵	85.3 ¹⁹ 90.8 ¹³	85.7 ¹⁹ 91.0 ¹³
SVHN (DenseNet)	CIFAR-10	90.689.5	97.797.8	99.199.2
	CIFAR-100	91.888.4	98.097.7	99.299.1
	Chars	72.370.5	95.294.5	98.598.3
	TinyImageNet	99.598.1	99.699.3	99.599.8
	LSUN	99.997.3	99.899.1	99.999.7
	Places365	94.391.9	98.398.2	98.199.3
	Textures	95.377.1	98.899.3	99.699.8
	Gaussian	100100	10099.9	100100
	iSUN	99.998.2	99.899.3	99.999.8
	average	93.7 ^{8.0} 92.3 ^{9.0}	98.6 ^{1.0} 98.3 ^{2.0}	99.3 ^{1.0} 99.4 ⁰
CIFAR-10 (ResNet)	CIFAR-100	36.521.5	84.563.3	84.358.1
	Chars	82.090.9	96.998.3	97.798.7
	TinyImageNet	96.294.3	99.298.0	99.296.7
	LSUN	98.297.7	99.599.2	99.598.9
	Places365	34.815.9	85.060.1	84.224.4
	SVHN	81.098.2	96.699.3	93.797.5
	Textures	81.781.6	96.793.4	98.294.3
	Gaussian	100100	100100	100100
	iSUN	96.895.3	99.398.6	99.398.1
	average	78.6 ²⁴ 77.3 ³²	95.3 ^{6.0} 90.0 ¹⁵	95.1 ^{6.0} 85.2 ²⁵
CIFAR-100 (ResNet)	CIFAR-10	3.05.0	61.059.6	63.760.6
	Chars	39.95.1	85.690.4	88.192.5
	TinyImageNet	88.786.2	97.697.3	97.697.3
	LSUN	91.388.6	98.097.8	98.398.0
	Places365	8.816	67.963.0	66.861.7
	SVHN	31.675.2	82.995.8	68.892.7
	Textures	65.978.1	91.995.6	95.297.6
	Gaussian	100100	100100	100100
	iSUN	87.989.4	97.497.8	97.697.7
	average	57.43865.1 ³³	86.9 ¹³ 88.6 ¹⁵	86.2 ¹⁴ 88.7 ¹⁵
SVHN (ResNet)	CIFAR-10	97.196.7	99.199.2	99.799.7
	CIFAR-100	97.596.2	99.199.1	99.799.6
	Chars	75.455.1	95.389.1	98.596.0
	TinyImageNet	99.999.6	99.999.9	99.999.9
	LSUN	10099.8	99.999.9	100100
	Places365	98.197.0	99.299.2	99.299.0
	Textures	98.998.4	99.699.6	99.999.9
	Gaussian	100100	99.9100	100100
	iSUN	10099.8	99.899.9	99.9100
	average	96.3 ^{8.0} 93.6 ¹⁴	99.1 ^{1.0} 98.4 ^{3.0}	99.6 ⁰ 99.3 ^{1.0}
Avg. and std. of avg. values		77.6 ⁷ 79.5 ¹⁰	92.8 ^{5.4} 93.4 ^{3.9}	92.3 ^{5.9} 92.9 ^{5.2}



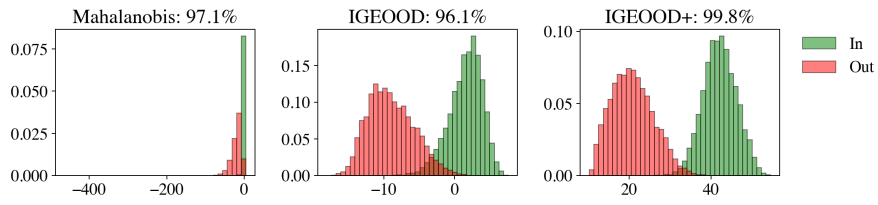
(a) DenseNet on CIFAR-10.



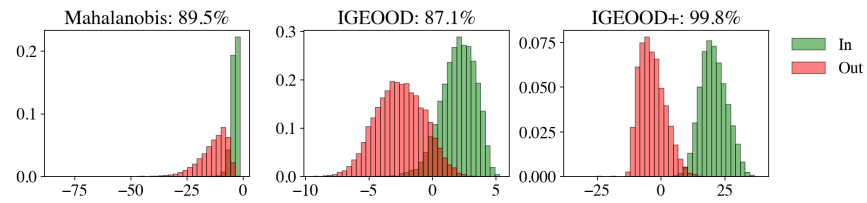
(b) DenseNet on CIFAR-100.



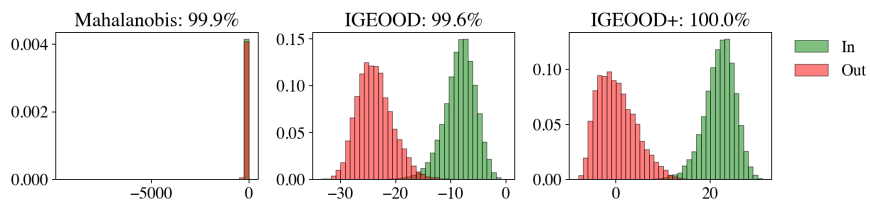
(c) DenseNet on SVHN.



(d) ResNet on CIFAR-10.



(e) ResNet on CIFAR-100.



(f) ResNet on SVHN.

Figure 2: WHITE-BOX setup with adversarial data validation. TinyImageNet as OOD dataset.

