

# Task Agnostic and Post-hoc Unseen Distribution Detection

Radhika Dua<sup>1</sup>   Seongjun Yang<sup>1</sup>   Yixuan Li<sup>2</sup>   Edward Choi<sup>1</sup>

<sup>1</sup>KAIST   <sup>2</sup>University of Wisconsin-Madison  
{radhikadua, seongjunyang, edwardchoi}@kaist.ac.kr  
sharonli@cs.wisc.edu

## Abstract

*Despite the recent advances in out-of-distribution(OOD) detection, anomaly detection, and uncertainty estimation tasks, there do not exist a task-agnostic and post-hoc approach. To address this limitation, we design a novel clustering-based ensembling method, called Task Agnostic and Post-hoc Unseen Distribution Detection (TAPUDD) that utilizes the features extracted from the model trained on a specific task. Explicitly, it comprises of TAP-Mahalanobis, which clusters the training datasets’ features and determines the minimum Mahalanobis distance of the test sample from all clusters. Further, we propose the Ensembling module that aggregates the computation of iterative TAP-Mahalanobis for a different number of clusters to provide reliable and efficient cluster computation. Through extensive experiments on real-world datasets, we observe that our task-agnostic approach can detect unseen samples effectively across diverse tasks and performs better or on-par with the existing task-specific baselines.*

## 1. Introduction

Deep neural networks have achieved phenomenal performance in diverse domains such as computer vision and healthcare [3,10,32]. However, they struggle to handle samples from an unseen distribution, leading to unreliable predictions and fatal errors in safety-critical applications. In an ideal situation, a robust model should be capable of making predictions on samples from the learned distributions, and at the same time, flag unknown inputs from unfamiliar distributions so that humans can make a responsible decision. For instance, in safety-critical tasks such as cancer detection, the machine learning assistant must issue a warning and hand over the control to the doctors when it detects an unusual sample that it has never seen during training. Thus, in practice, it is important for a model to know when *not* to predict. This task of detecting samples from an unseen distribution is referred to as out-of-distribution (OOD) detection [5,17,18,21,24,26,27,31,41,47].

Most of these OOD detection methods mainly focusing

on classification tasks have shown great success. However, they are not directly applicable to other tasks like regression. Although a few bayesian and non-bayesian techniques [11,14,23,30] estimate uncertainty in regression tasks, they are not post-hoc as it often requires a modification to the training pipeline, or multiple trained copies of the model, or training a model with an optimal dropout rate. This raises an under-explored question:

*Can we design a task-agnostic, and post-hoc approach for unseen distribution detection ?*

Motivated by this, we propose a novel clustering-based ensembling framework, “Task Agnostic and Post-hoc Unseen Distribution Detection (TAPUDD)”, which comprises of two modules, *TAP-Mahalanobis* and *Ensembling*. *TAP-Mahalanobis* partitions the training datasets’ features into clusters and then determines the minimum Mahalanobis distance of a test sample from all the clusters. The *Ensembling* module aggregates the outputs obtained from *TAP-Mahalanobis* iteratively for a different number of clusters. It enhances reliability and eliminates the need to determine an optimal number of clusters. As TAPUDD is a post-hoc approach and doesn’t require training the model, it is more efficient and easy to deploy in real-world. We demonstrate the efficacy of our approach by extensively evaluating it on synthetic and real-world datasets for diverse tasks.

## 2. Related Work

**Out-of-distribution Detection.** Recent works have introduced reconstruction-error based [7, 8, 37, 38, 40, 48], density-based [6,9,13,31,34,39,42], and self-supervised [1, 12, 19, 41] OOD detection methods. Other efforts include post-hoc methods [5,17,18,24,26,27,33] that do not require modification to the training procedure. However, there is no approach that is post-hoc and does not require the class label information of the training data.

**Uncertainty Estimation.** Research in this direction primarily estimates the uncertainty to enhance the robustness of networks in regression tasks. Well-known methods to estimate uncertainty include bayesian [2, 14, 16, 22, 25, 28–

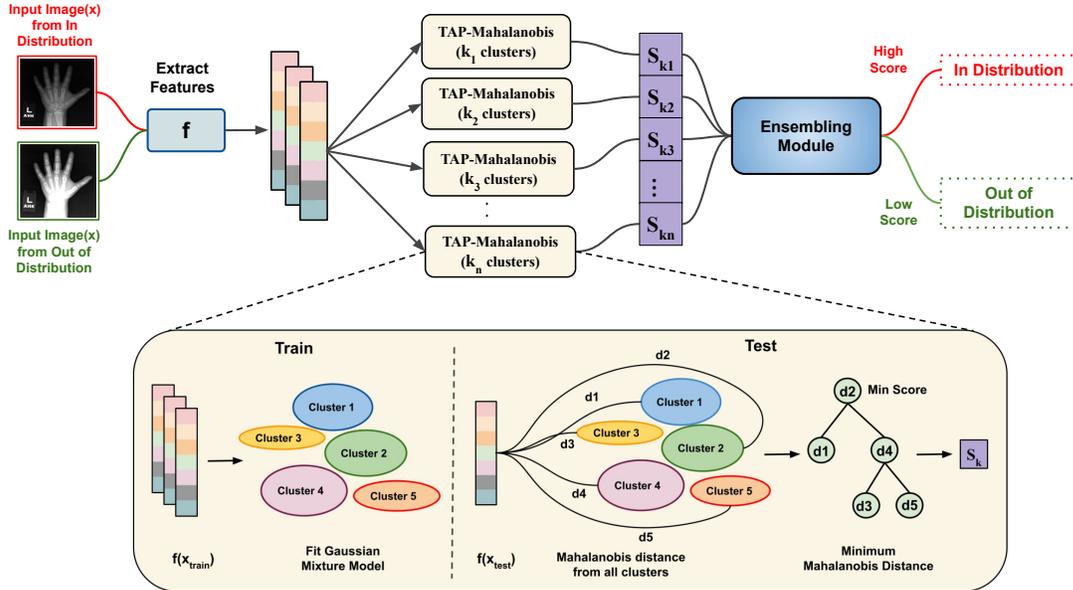


Figure 1. **TAPUDD**. Our method first extracts the features of an input image  $\mathbf{x}$  from the feature extractor  $f$  of a model trained on a specific task. *TAP-Mahalanobis* module then uses the extracted features  $f(\mathbf{x}_{train})$  to fit the gaussian mixture model and computes the minimum mahalanobis distance  $S_k$  for the given feature vector  $f(\mathbf{x}_{test})$ . Further, the *Ensembling* module aggregates the mahalanobis distance ( $S_{k_1}$  to  $S_{k_n}$ ) obtained from iterative computation of *TAP-Mahalanobis* for different number of clusters ( $k_1$  to  $k_n$ ) to enhance the reliability.

30, 35, 43, 46] and non-bayesian [11, 23] approaches, which have shown remarkable success. However, they require significant modification to the training pipeline, multiple trained copies of the model, and are not post-hoc.

**Anomaly Detection.** This task aims to detect anomalous samples shifted from the defined normality. Prior work [4, 8, 12, 37, 40, 45, 48] proposed methods to solve anomaly detection. However, more recently, [1, 20, 41, 44] proposed a unified method to solve both OOD detection and anomaly detection. Nonetheless, these methods require end-to-end training and are not post-hoc.

There exist no unified approach to enhance the reliability of neural networks across distinct tasks like classification, regression, etc. In contrast to all the aforementioned efforts, our work presents a post-hoc, and task-agnostic approach to detect unknown samples across varied tasks.

### 3. TAPUDD: Task Agnostic and Post-hoc Unseen Distribution Detection

We propose a novel, **Task Agnostic and Post-hoc Unseen Distribution Detection (TAPUDD)** method, as shown in Fig. 1. The method comprises of two main modules *TAP-Mahalanobis* and *Ensembling*.

**TAP-Mahalanobis.** Given training samples  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ , we extract the features of the in-distribution data from a model trained for a specific task using a feature extractor  $f$ . We then pass these features to the *TAP-Mahalanobis* module. It first partition the features of the

in-distribution data into  $K$  clusters using Gaussian Mixture Model (GMM) with “full” covariance. Then, we model the features in each cluster independently as multivariate gaussian and compute the empirical cluster mean and covariance of training samples  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  and their corresponding cluster labels  $C = \{c_1, \dots, c_N\}$  as:

$$\mu_c = \frac{1}{N_c} \sum_{i:c_i=c} f(\mathbf{x}_i), \Sigma_c = \frac{1}{N_c} \sum_{i:c_i=c} (f(\mathbf{x}_i) - \mu_c)(f(\mathbf{x}_i) - \mu_c)^T,$$

where  $f(\mathbf{x}_i)$  denotes the penultimate layer features of an input sample  $\mathbf{x}_i$  from a cluster  $c_i$ .

Then, given a test sample,  $\mathbf{x}_{test}$ , we obtain the negative of the minimum of the Mahalanobis distance from the center of the clusters as follows:

$$\mathcal{S}_{TAP-Mahalanobis} = -\min_c (f(\mathbf{x}_{test}) - \mu_c)^T \Sigma_c^{-1} (f(\mathbf{x}_{test}) - \mu_c),$$

where  $f(\mathbf{x}_{test})$  denotes the penultimate layer features of a test sample  $\mathbf{x}_{test}$ . We then use the score  $\mathcal{S}_{TAP-Mahalanobis}$  to distinguish between ID and OOD samples. To align with the conventional notion of having high score for ID samples and low score for OOD samples, negative sign is applied.

However, it is not straightforward to determine the number of clusters  $K$  for which the OOD detection performance of *TAP-Mahalanobis* is optimal for different tasks and datasets. Therefore, we present an *Ensembling* module.

**Ensembling.** This module not only eliminates the need to determine the optimal value of  $K$  but also provides more reliable results. We obtain *TAP-Mahalanobis* scores for different values of  $K \in [k_1, k_2, k_3, \dots, k_n]$  and average them to obtain an ensemble score,  $\mathcal{S}_{Ensemble}$ . This ensures that a sample is detected as OOD only if a majority of the participants in ensembling agrees with each other.

Brightness	Baselines							Ours (Task-Agnostic)	
	MSP [18]	ODIN [26]	Energy [27]	MB [24]	KL [17]	MOS [21] (K = 8)	Gram [5]	TAP-MB (K = 8)	TAPUDD (Average)
0.0	88.7±4.8	88.7±4.8	88.2±5.3	99.9±0.1	26.3±32.8	89.3±5.5	99.3±1.4	99.9±0.1	100.0±0.1
0.2	66.1±3.5	66.1±3.5	66.0±3.7	87.5±4.5	44.5±3	65.9±3.2	61.0±3.3	86.8±4.7	87.3±5.2
0.4	56.3±1.4	56.4±1.4	56.2±1.7	70.5±3.8	46.9±1.2	56.4±1.1	53.4±1.1	69.6±3.7	70.1±4.5
0.6	52.4±0.8	52.4±0.8	52.3±0.9	59.9±2.5	48.2±1	52.5±0.8	51.4±0.5	59.3±2.5	59.4±2.7
0.8	50.4±0.4	50.4±0.4	50.4±0.4	52.2±1.4	48.8±0.6	50.5±0.3	50.2±0.5	52.0±1.7	52.0±1.6
1.0	50.0±0.0	50.0±0.0	50.0±0.0	50.0±0.0	50.0±0.0	50.0±0.0	50.0±0.0	50.0±0.0	50.0±0.0
1.2	51.7±0.4	51.7±0.4	51.7±0.4	55.4±1.6	49.2±0.5	51.7±0.5	51.1±0.6	56.1±1.5	56.0±1.5
1.4	55.8±0.8	55.8±0.8	55.8±0.8	62.9±2.1	48.2±1.2	55.8±0.8	53.6±1.1	63.7±2.0	63.5±2.1
1.6	59.7±1.3	59.7±1.3	59.8±1.4	70.2±2.7	47.5±1.7	59.6±1.1	55.9±1.2	70.9±2.8	70.7±2.9
1.8	63.1±2.0	63.1±2.1	63.2±2.2	76.5±2.9	48.3±2.4	62.8±1.7	58.1±1.5	76.9±3.4	76.6±3.5
2.0	65.5±3.2	65.6±3.2	65.7±3.5	81.6±2.7	49.8±2.9	65.1±2.6	60.5±1.8	81.8±3.7	81.4±3.8
2.5	69.5±6.5	69.5±6.5	69.6±6.8	90.4±2.5	51.6±4.9	69.0±5.5	65.4±4.4	89.9±3.8	89.6±4.1
3.0	72.5±8.7	72.5±8.7	72.6±9.0	94.8±1.8	51.3±5.4	72.0±7.6	69.6±5.9	93.9±3.8	93.6±4.0
3.5	73.7±9.7	73.7±9.7	73.6±10	96.8±1.3	52.0±6.1	73.0±8.8	72.2±6.8	95.5±3.8	95.4±3.7
4.0	75.8±9.5	75.8±9.5	75.7±9.8	97.8±0.8	50.5±7.2	75.3±8.8	75.1±7.9	96.5±3.6	96.5±3.2
4.5	78.1±7.9	78.1±7.9	78.0±8.3	98.5±0.5	47.4±8.2	77.8±7.5	78.4±7.1	97.3±3.0	97.4±2.4
5.0	79.9±6.4	79.9±6.4	79.8±6.9	98.8±0.4	44.9±8.4	79.8±6.1	80.4±6.6	97.9±2.5	98.0±1.7
5.5	81.4±5.6	81.4±5.6	81.3±6.2	99.0±0.4	44.1±8.7	81.3±5.4	82.4±6.6	98.2±2.2	98.4±1.2
6.0	82.5±5.1	82.5±5.1	82.5±5.6	99.1±0.4	43.6±8.6	82.4±4.9	83.9±6.3	98.5±1.9	98.7±0.9
6.5	83.2±4.9	83.2±4.9	83.2±5.4	99.2±0.4	44.3±8.2	83.1±4.6	85.0±6.2	98.7±1.7	98.9±0.7
Average	67.8	67.8	67.8	82.1	46.9	67.7	66.8	<b>81.7</b>	<b>82.0</b>

Table 1. NAS detection performance in binary classification task for NAS shift of brightness in RSNA boneage dataset measured by AU-ROC. Highlighted row presents the performance on ID data. MB and TAP-MB refers to Mahalanobis and *TAP-Mahalanobis*, respectively. Our task-agnostic approach significantly outperforms all baselines (except MB) and is comparable to MB. Note that **MB is task-specific** and cannot be used in tasks other than classification.

**Remark.** GMM is more flexible in learning the cluster shape in contrast to K-means, which learns spherical cluster shapes. Consequently, K-means performs poorly when detecting OOD samples near the cluster. Other popular clustering methods such as agglomerative clustering or DBSCAN are less compatible with Mahalanobis distance and require careful hyperparameter adjustment, such as the linking strategies for agglomerative clustering or the epsilon value for DBSCAN.

## 4. Experiments and Results

In this section, we validate TAPUDD by conducting experiments on 2-D synthetic dataset (Sec. 4.1). To further bolster the effectiveness of our method, we present empirical evidence to validate TAPUDD on several real-world tasks, including binary classification (Sec. 4.2), and regression (Sec. 4.3). For real-world tasks, we evaluate on Natural Attribute-based Shift (NAS) detection dataset [36]. In NAS detection, a sample is shifted from the training distribution based on attributes like brightness, age, etc.

### 4.1. Evaluation on Synthetic Datasets

**Experimental Details.** We generate synthetic datasets in  $\mathbb{R}^2$  for multi-class classification task. The in-distribution (ID) data  $\mathbf{x} \in \mathcal{X} = \mathbb{R}^2$  is sampled from a Gaussian mixture model. All the samples except the ID samples in the 2-D plane represent the OOD samples. We consider the 2-D sample as the penultimate layer features on which we can directly apply OOD detection methods like TAPUDD.

**TAPUDD outperforms TAP-Mahalanobis.** We present

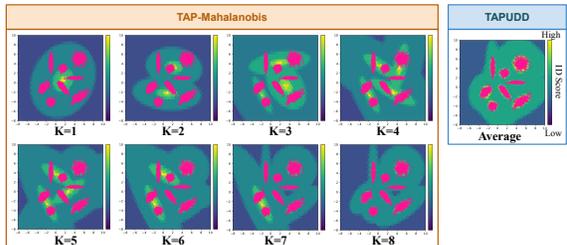


Figure 2. ID score landscape of *TAP-Mahalanobis* for different values of  $K$  (i.e., number of clusters); and TAPUDD on synthetic 2D multi-class classification dataset. A sample is deemed as OOD when it has a **low ID score**. The **Pink Points** represent the in-distribution data. Results demonstrate that *TAP-Mahalanobis* does not perform well for some values of  $K$  whereas TAPUDD perform better or on-par with *TAP-Mahalanobis*.

a comparison to demonstrate the effectiveness of TAPUDD against *TAP-Mahalanobis* in Fig. 2. We present the ID score landscape of *TAP-Mahalanobis* for different values of  $K$  and TAPUDD for multi-class classification in a 2-D synthetic dataset. The **Pink Points** represent the ID data. We observe that for certain values of  $K$ , *TAP-Mahalanobis* fails to detect some OOD samples. However, TAPUDD effectively detect OOD samples and performs better, or on par, with *TAP-Mahalanobis*. Thus, TAPUDD eliminates the necessity of choosing the optimal value of  $K$ .

### 4.2. NAS Detection in Binary Classification

**Experimental Details.** We use the RSNA Bone Age dataset [15], composed of left-hand X-ray images of the patient and their gender and age (0 to 20 years). We alter the brightness of the X-ray images by a factor between 0 and

6.5 and form 20 different NAS datasets to reflect the X-ray imaging set-ups in different hospitals following [36]. In-distribution (ID) data consists of images with a brightness factor 1.0. We trained a ResNet18 model using the cross-entropy loss and assessed it on the ID test set composed of images with a brightness factor of 1.0. Further, we evaluate the NAS detection performance of our method and compare it with representative task-specific OOD detection methods on NAS datasets. For NAS detection, we measure the area under the receiver operating characteristic curve (AUROC), a commonly used metric for OOD detection.

**Results.** The ID classification accuracy averaged across 10 seeds of the gender classifier trained using cross-entropy loss is 91.60. We compare the NAS detection performance of our proposed approach with competitive post-hoc OOD detection methods in literature in Tab. 1. As expected, the NAS detection performance of our approach and all baselines except KL Matching increase as the shift in the brightness attribute increases. We also observe that our approaches, *TAPUDD* and *TAP-Mahalanobis* are more sensitive to NAS samples compared to competitive baselines, including Maximum Softmax Probability [18], ODIN [26], Mahalanobis distance [24], energy score [27], Gram matrices [5], MOS [21], and KL matching [17]. All these task-specific baselines require the label information of the training dataset for OOD detection and cannot be used directly in tasks other than classification. In contrast, our proposed task-agnostic approach does not require the access to class label information and it can be used across different tasks.

### 4.3. NAS Detection in Regression

**Experimental Details.** We use the RSNA Bone Age dataset (described in Sec. 4.2) and solve the age prediction task. In this task, the objective is to automatically predict the patient’s age given a hand X-ray image as an input. As described in Sec. 4.2, we vary the brightness and form 20 different NAS datasets. In-distribution (ID) data comprises images of brightness factor 1.0 (unmodified images). We train a ResNet18 with MSE loss and evaluate it on the test set composed of images with a brightness factor 1.0. Further, we evaluate the NAS detection performance of our proposed method and compare its performance with representative bayesian and non-bayesian uncertainty estimation methods on NAS datasets with attribute shift of brightness.

**Results.** The in-distribution Mean Absolute Error (MAE) in year averaged across 10 seeds of the Resnet18 model trained using MSE loss is 0.801. We compare the NAS detection performance of our proposed approach with well-known uncertainty estimation methods, namely Deep Ensemble (DE) [23], Monte Carlo Dropout (MCD) [11], and SWAG [30]. Although DE, MCD, and SWAG are not applicable to a pre-trained model, we compare against these baselines as a benchmark, as it has shown strong OOD de-

Brightness	Baselines			Ours (Task-Agnostic)	
	DE [23]	MCD [11]	SWAG* [30]	TAP-MB (K = 8)	TAPUDD (Average)
0.0	100.0±NA	6.9±NA	99.9±NA	100.0±0.1	100.0±0.0
0.2	57.0±NA	45.5±NA	51.4±NA	87.9±6.1	88.8±6.7
0.4	51.3±NA	50.8±NA	49.8±NA	64.5±6.9	66.6±5.0
0.6	50.7±NA	49.7±NA	49.5±NA	54.6±4.4	55.1±2.5
0.8	50.5±NA	49.9±NA	49.7±NA	48.9±1.7	49.2±1.0
1.0	50.0±NA	49.8±NA	50.0±NA	50.0±0.0	50.0±0.0
1.2	50.3±NA	48.5±NA	50.8±NA	57.6±1.8	57.8±1.9
1.4	54.5±NA	46.7±NA	55.8±NA	68.4±3.4	68.4±3.4
1.6	58.6±NA	44.5±NA	63.5±NA	78.7±3.6	78.6±3.7
1.8	64.9±NA	41.6±NA	71.6±NA	86.4±3.5	86.3±3.6
2.0	75.8±NA	38.4±NA	79.3±NA	91.9±3.0	91.7±3.2
2.5	95.6±NA	31.1±NA	89.8±NA	97.5±1.5	97.4±1.4
3.0	98.4±NA	25.8±NA	90.7±NA	99.0±0.6	99.0±0.5
3.5	99.3±NA	21.7±NA	93.7±NA	99.4±0.3	99.4±0.3
4.0	99.8±NA	18.0±NA	96.4±NA	99.6±0.3	99.6±0.2
4.5	100.0±NA	14.9±NA	97.4±NA	99.7±0.2	99.7±0.1
5.0	100.0±NA	11.7±NA	98.1±NA	99.8±0.1	99.7±0.1
5.5	100.0±NA	9.7±NA	98.5±NA	99.8±0.1	99.8±0.2
6.0	100.0±NA	7.9±NA	98.7±NA	99.8±0.1	99.8±0.2
6.5	100.0±NA	7.0±NA	98.9±NA	99.8±0.2	99.8±0.3
Average	77.8	31.0	76.7	<b>84.2</b>	<b>84.3</b>

Table 2. NAS detection performance in regression task (age prediction) for NAS shift of brightness in RSNA boneage dataset measured by AUROC. Highlighted row presents the performance on the ID dataset. DE, MCD, TAP-MB, and NA denotes Deep Ensemble, Monte Carlo Dropout, *TAP-Mahalanobis*, and Not Applicable respectively. SWAG\* = SWAG + Deep Ensemble.

tection performance across regression examples. For DE, we retrain 10 models of the same architecture using MSE loss from different initializations. Since SWAG is not directly applicable for OOD detection, we apply SWAG\* which is a combination of deep ensembling on top of SWAG. From Tab. 2, as expected, we observe that the NAS detection performance of our approach and all baselines increase as the shift in the brightness attribute increases. We also observe that our proposed approaches, *TAPUDD* and *TAP-Mahalanobis*, are more sensitive to NAS samples and effectively detect them compared to the baselines.

## 5. Conclusion

In this work, we propose a task-agnostic and post-hoc approach, *TAPUDD*, to detect samples from the unseen distribution. *TAPUDD* is a clustering-based ensembling approach composed of *TAP-Mahalanobis* and *Ensembling* modules. *TAP-Mahalanobis* module groups the semantically similar training samples into clusters and determines the minimum Mahalanobis distance of the test sample from the clusters. To enhance reliability and to eliminate the necessity to determine the optimal number of clusters for *TAP-Mahalanobis*, the *Ensembling* module aggregates the distances obtained from the *TAP-Mahalanobis* module for different values of clusters. We validate the effectiveness of our approach by conducting extensive experiments on diverse datasets and tasks. As future work, it would be interesting to extensively evaluate *TAPUDD* to detect unseen samples in text, 3D vision, and healthcare.

## References

- [1] Liron Bergman and Yedid Hoshen. Classification-based anomaly detection for general data. *ArXiv*, abs/2005.02359, 2020. [1](#), [2](#)
- [2] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural networks. *ArXiv*, abs/1505.05424, 2015. [1](#)
- [3] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Praseon Goyal, Lawrence D. Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, Xin Zhang, Jake Zhao, and Karol Zieba. End to end learning for self-driving cars, 2016. [1](#)
- [4] Raghavendra Chalapathy, Aditya Krishna Menon, and Sanjay Chawla. Anomaly detection using one-class neural networks. *ArXiv*, abs/1802.06360, 2018. [2](#)
- [5] Sastry Shama Chandramouli and Oore Sageev. Detecting out-of-distribution examples with gram matrices. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2020. [1](#), [3](#), [4](#)
- [6] Hyunsun Choi, Eric Jang, and Alexander A Alemi. Waic, but why? generative ensembles for robust anomaly detection. *arXiv preprint arXiv:1810.01392*, 2018. [1](#)
- [7] Sung-Ik Choi and Sae-Young Chung. Novelty detection via blurring. *ArXiv*, abs/1911.11943, 2020. [1](#)
- [8] Lucas Deecke, Robert A. Vandermeulen, Lukas Ruff, Stephan Mandt, and M. Kloft. Image anomaly detection with generative adversarial networks. In *Proceedings of the European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML-PKDD)*, 2018. [1](#), [2](#)
- [9] Yilun Du and Igor Mordatch. Implicit generation and modeling with energy based models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. [1](#)
- [10] Andre Esteva, Brett Kuprel, Roberto Novoa, Justin Ko, Swetter Susan, Helen Balu, M, and Sebastian Thrun. Dermatologist-level classification of skin cancer with deep neural networks. 2016. [1](#)
- [11] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2016. [1](#), [2](#), [4](#)
- [12] Izhak Golan and Ran El-Yaniv. Deep anomaly detection using geometric transformations. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018. [1](#), [2](#)
- [13] Will Grathwohl, Kuan-Chieh Wang, Jörn-Henrik Jacobsen, David Kristjanson Duvenaud, Mohammad Norouzi, and Kevin Swersky. Your classifier is secretly an energy based model and you should treat it like one. *ArXiv*, abs/1912.03263, 2020. [1](#)
- [14] Alex Graves. Practical variational inference for neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2011. [1](#)
- [15] Safwan S. Halabi, Luciano M. Prevedello, Jayashree Kalpathy-Cramer, Artem B. Mamonov, Alexander Bilbily, Mark Cicero, Ian Pan, Lucas Araújo Pereira, Rafael Teixeira Sousa, Nitamar Abdala, Felipe Campos Kitamura, Hans H. Thodberg, Leon Chen, George Shih, Katherine Andriole, Marc D. Kohli, Bradley J. Erickson, and Adam E. Flanders. The rsna pediatric bone age machine learning challenge. *Radiology*, 290(3):498–503, 2019. [3](#)
- [16] Leonard Hasenclever, Stefan Webb, Thibaut Lienart, Sebastian J. Vollmer, Balaji Lakshminarayanan, Charles Blundell, and Yee Whye Teh. Distributed bayesian learning with stochastic natural gradient expectation propagation and the posterior server. *ArXiv*, abs/1512.09327, 2017. [1](#)
- [17] Dan Hendrycks, Steven Basart, Mantas Mazeika, Mohammadreza Mostajabi, Jacob Steinhardt, and Dawn Xiaodong Song. A benchmark for anomaly segmentation. *ArXiv*, abs/1911.11132, 2019. [1](#), [3](#), [4](#)
- [18] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017. [1](#), [3](#), [4](#)
- [19] Dan Hendrycks, Mantas Mazeika, Saurav Kadavath, and Dawn Xiaodong Song. Using self-supervised learning can improve model robustness and uncertainty. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. [1](#)
- [20] José Miguel Hernández-Lobato and Ryan P. Adams. Probabilistic backpropagation for scalable learning of bayesian neural networks. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2015. [2](#)
- [21] Rui Huang and Yixuan Li. Mos: Towards scaling out-of-distribution detection for large semantic space. *Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8706–8715, 2021. [1](#), [3](#), [4](#)
- [22] Anoop Korattikara, Vivek Rathod, Kevin Murphy, and Max Welling. Bayesian dark knowledge, 2015. [1](#)
- [23] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. [1](#), [2](#), [4](#)
- [24] Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018. [1](#), [3](#), [4](#)
- [25] Yingzhen Li, José Miguel Hernández-Lobato, and Richard E. Turner. Stochastic expectation propagation. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2015. [1](#)
- [26] Shiyu Liang, Yixuan Li, and Rayadurgam Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018. [1](#), [3](#), [4](#)
- [27] Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection. *Advances in Neural Information Processing Systems (NeurIPS)*, 33, 2020. [1](#), [3](#), [4](#)
- [28] Christos Louizos and Max Welling. Structured and efficient variational deep learning with matrix gaussian posteriors. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2016. [1](#)
- [29] David J. C. Mackay. Bayesian methods for adaptive models. 1992. [1](#)

- [30] Wesley J Maddox, Pavel Izmailov, Timur Garipov, Dmitry P Vetrov, and Andrew Gordon Wilson. A simple baseline for bayesian uncertainty in deep learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 13153–13164, 2019. 1, 4
- [31] Ahsan Mahmood, Junier Oliva, and Martin Andreas Styner. Multiscale score matching for out-of-distribution detection. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020. 1
- [32] Tomás Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2013. 1
- [33] Eric T. Nalisnick, Akihiro Matsukawa, Yee Whye Teh, Dilan Görür, and Balaji Lakshminarayanan. Do deep generative models know what they don’t know? *ArXiv*, abs/1810.09136, 2019. 1
- [34] Eric T. Nalisnick, Akihiro Matsukawa, Yee Whye Teh, and Balaji Lakshminarayanan. Detecting out-of-distribution inputs to deep generative models using a test for typicality. *ArXiv*, abs/1906.02994, 2019. 1
- [35] Radford M. Neal. Bayesian learning for neural networks. 1995. 1
- [36] Jeonghoon Park, Jimin Hong, Radhika Dua, Daehoon Gwak, Yixuan Li, Jaegul Choo, and E. Choi. Natural attribute-based shift detection. *ArXiv*, abs/2110.09276, 2021. 3, 4
- [37] Pramuditha Perera, Ramesh Nallapati, and Bing Xiang. Ocgan: One-class novelty detection using gans with constrained latent representations. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (cvpr)*, pages 2893–2901, 2019. 1, 2
- [38] Stanislav Pidhorskyi, Ranya Almohsen, Donald A. Adjeroh, and Gianfranco Doretto. Generative probabilistic novelty detection with adversarial autoencoders. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018. 1
- [39] Jie Ren, Peter J Liu, Emily Fertig, Jasper Snoek, Ryan Poplin, Mark A DePristo, Joshua V Dillon, and Balaji Lakshminarayanan. Likelihood ratios for out-of-distribution detection. *arXiv preprint arXiv:1906.02845*, 2019. 1
- [40] 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 *International Conference on Information Processing in Medical Imaging*, 2017. 1, 2
- [41] Vikash Sehwal, Mung Chiang, and Prateek Mittal. {SSD}: A unified framework for self-supervised outlier detection. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021. 1, 2
- [42] Joan Serrà, David Álvarez, V. Gómez, Olga Slizovskaia, José F. Núñez, and Jordi Luque. Input complexity and out-of-distribution detection with likelihood-based generative models. *ArXiv*, abs/1909.11480, 2020. 1
- [43] Jost Tobias Springenberg, Aaron Klein, Stefan Falkner, and Frank Hutter. Bayesian optimization with robust bayesian neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2016. 1
- [44] Jihoon Tack, Sangwoo Mo, Jongheon Jeong, and Jinwoo Shin. Csi: Novelty detection via contrastive learning on distributionally shifted instances. *ArXiv*, abs/2007.08176, 2020. 2
- [45] Siqi Wang, Yijie Zeng, Xinwang Liu, En Zhu, Jianping Yin, Chuanfu Xu, and M. Kloft. Effective end-to-end unsupervised outlier detection via inlier priority of discriminative network. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 2
- [46] Max Welling and Yee Whye Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2011. 1
- [47] Zhisheng Xiao, Qing Yan, and Yali Amit. Do we really need to learn representations from in-domain data for outlier detection? *ArXiv*, abs/2105.09270, 2021. 1
- [48] Bo Zong, Qi Song, Martin Renqiang Min, Wei Cheng, Cristian Lumezanu, Dae ki Cho, and Haifeng Chen. Deep autoencoding gaussian mixture model for unsupervised anomaly detection. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018. 1, 2