

Generalized Out-of-Distribution Detection: A Survey

Jingkang Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu

Abstract—Out-of-distribution (OOD) detection is critical to ensuring the reliability and safety of machine learning systems. For instance, in autonomous driving, we would like the driving system to issue an alert and hand over the control to humans when it detects unusual scenes or objects that it has never seen during training time and cannot make a safe decision. The term, OOD detection, first emerged in 2017 and since then has received increasing attention from the research community, leading to a plethora of methods developed, ranging from classification-based to density-based to distance-based ones. Meanwhile, several other problems, including anomaly detection (AD), novelty detection (ND), open set recognition (OSR), and outlier detection (OD), are closely related to OOD detection in terms of motivation and methodology. Despite common goals, these topics develop in isolation, and their subtle differences in definition and problem setting often confuse readers and practitioners. In this survey, we first present a unified framework called *generalized OOD detection*, which encompasses the five aforementioned problems, *i.e.*, AD, ND, OSR, OOD detection, and OD. Under our framework, these five problems can be seen as special cases or sub-tasks, and are easier to distinguish. We then review each of these five areas by summarizing their recent technical developments, with a special focus on OOD detection methodologies. We conclude this survey with open challenges and potential research directions.

Index Terms—Out-of-Distribution Detection, Open Set Recognition, Anomaly Detection, Novelty Detection, Outlier Detection

1 INTRODUCTION

A trustworthy visual recognition system should not only produce accurate predictions on known context, but also detect unknown examples and reject them (or hand them over to human users for safe handling) [1], [2], [3], [4], [5], [6], [7], [8]. For instance, a well-trained food classifier should be able to detect non-food images such as selfies uploaded by users, and reject such input instead of blindly classifying them into existing food categories. In safety-critical applications such as autonomous driving, the driving system must issue a warning and hand over the control to drivers when it detects unusual scenes or objects it has never seen during training.

Most existing machine learning models are trained based on the closed-world assumption [9], [10], where the test data is assumed to be drawn *i.i.d.* from the same distribution as the training data, known as in-distribution (ID). However, when models are deployed in an *open-world* scenario [11], test samples can be out-of-distribution (OOD) and therefore should be handled with caution. The distributional shifts can be caused by semantic shift (*e.g.*, OOD samples are drawn from different classes) [12], or covariate shift (*e.g.*, OOD samples from a different domain) [13], [14], [15].

The detection of semantic distribution shift (*e.g.*, due to the occurrence of new classes) is the focal point of OOD detection tasks, where the label space \mathcal{Y} can be different between ID and OOD data and hence the model should not

• J. Yang, K. Zhou, and Z. Liu are with S-Lab, School of Computer Science and Engineering, Nanyang Technological University, Singapore, 639798. E-mail: {jingkang001,kaiyang.zhou,ziwei.liu}@ntu.edu.sg
 • Y. Li is with Department of Computer Sciences, University of Wisconsin-Madison, Madison, WI, United States, 53706. E-mail: sharonli@cs.wisc.edu

Manuscript updated July 26, 2022. Discussions, comments, and questions are all welcomed in <https://github.com/jingkang50/OODSurvey/discussions>.

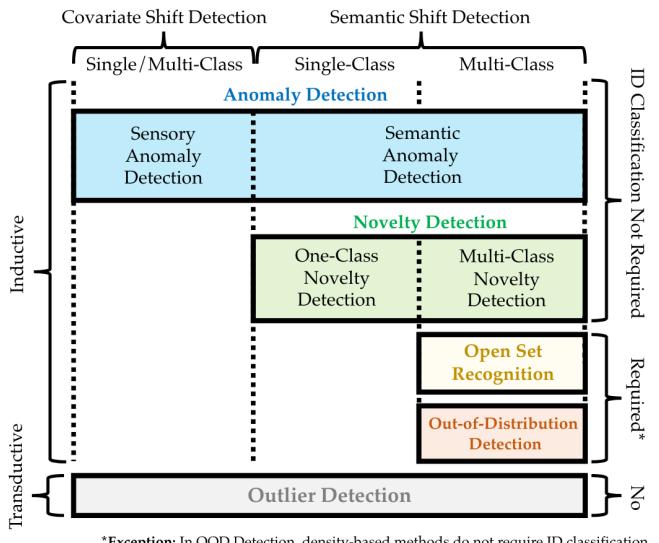


Fig. 1: Taxonomy of generalized OOD detection framework, illustrated by classification tasks. Four bases are used for the task taxonomy: **1)** Distribution shift to detect: the task focuses on detecting covariate shift or semantic shift; **2)** ID data type: the ID data contains one single class or multiple classes; **3)** Whether the task requires ID classification; **4)** Transductive learning task requires all observations; inductive tasks follow the train-test scheme. Note that ND is often interchangeable with AD, but ND is more concerned with semantic anomalies. OOD detection is generally interchangeable with OSR for classification tasks.

make any prediction. In addition to OOD detection, several problems adopt the “open-world” assumption and have a similar goal of identifying OOD examples. These include

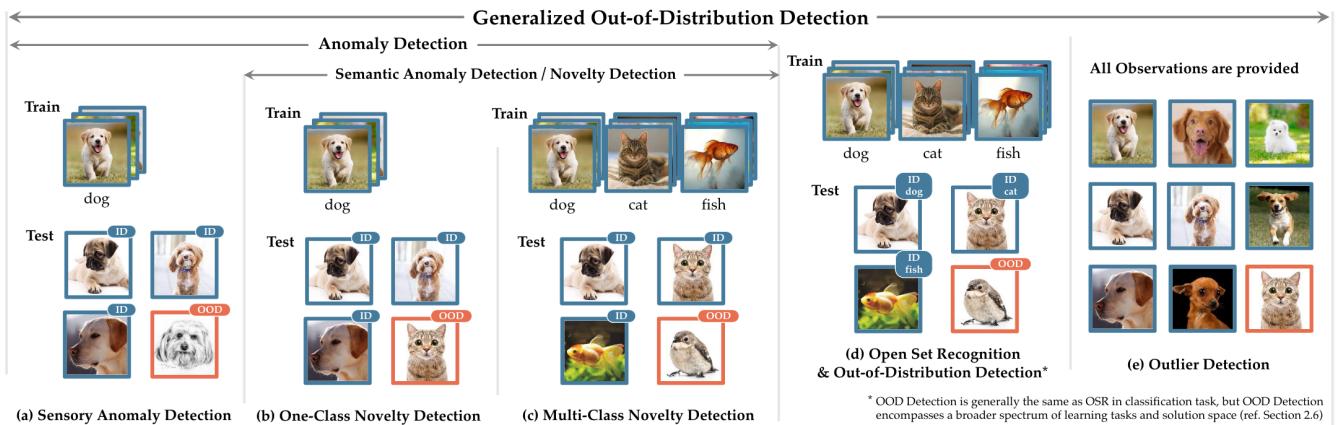


Fig. 2: Illustration of sub-tasks under generalized OOD detection framework with vision tasks. Tags on test images refer to model’s expected predictions. (a) In *sensory anomaly detection*, test images with covariate shift will be considered as OOD. No semantic shift occurs in this setting. (b) In *one-class novelty detection*, normal/ID images belong to one class. Test images with semantic shift will be considered as OOD. (c) In *multi-class novelty detection*, ID images belong to multiple classes. Test images with semantic shift will be considered as OOD. Note that (b) and (c) compose novelty detection, which is identical to the topic of semantic anomaly detection. (d) *Open set recognition* is identical to multi-class novelty detection in the task of detection, with the only difference that open set recognition further requires ID classification. *Out-of-distribution detection* solves the same problem as open-set recognition. It canonically aims to detect test samples with semantic shift without losing the ID classification accuracy. However, OOD Detection encompasses a broader spectrum of learning tasks and solution space. (e) *Outlier detection* does not follow a train-test scheme. All observations are provided. It fits in the generalized OOD detection framework by defining the majority distribution as ID. Outliers can have any distribution shift from the majority.

outlier detection (OD) [16], [17], [18], [19], anomaly detection (AD) [20], [21], [22], [23], novelty detection (ND) [24], [25], [26], [27], and open set recognition (OSR) [28], [29], [30]. While all these problems are related to each other by sharing similar motivations, subtle differences exist among the *sub-topics* in terms of the specific definition. However, the lack of a comprehensive understanding of the relationship between the different sub-topics leads to confusion for both researchers and practitioners. Even worse, these sub-topics, which are supposed to be compared and learned from each other, are developing in isolation.

In this survey, we for the first time clarify the similarities and differences between these problems, and present a unified framework termed *generalized OOD detection*. Under this framework, the five problems (*i.e.*, AD, ND, OSR, OOD detection, and OD) can be viewed as special cases or sub-topics. We further conduct a literature review for each sub-topic, with a special focus on the OOD detection task. To sum up, we make three contributions to the research community:

- 1) **A Unified Framework:** For the first time, we systematically review five closely related topics of AD, ND, OSR, OOD detection, and OD, and present a unified framework of *generalized OOD detection*. Under this framework, the similarities and differences of the five sub-topics can be systematically compared and analyzed. We hope our unification helps the community better understand these problems and correctly position their research in the literature.
- 2) **A Comprehensive Survey for OOD Detection:** Notice the existence of comprehensive surveys on AD, ND, OSR, and OD methodologies in recent years [20], [21], [22], [23], [29], this survey highlights

the overview of OOD detection methods and thus complements existing surveys. For completeness, methodologies of other sub-topics are also reviewed briefly. We hope our survey can help readers build a better understanding of the developments for each problem and their connections.

- 3) **Future Research Directions:** We draw readers’ attention to some problems or limitations that remain in the current generalized OOD detection field. We conclude this survey with discussions on open challenges and opportunities for future research.

2 GENERALIZED OOD DETECTION

Framework Overview In this section, we introduce a unified framework termed *generalized OOD detection*, which encapsulates five related sub-topics: anomaly detection (AD), novelty detection (ND), open set recognition (OSR), out-of-distribution (OOD) detection, and outlier detection (OD). These sub-topics can be similar in the sense that they all define a certain *in-distribution*, with the common goal of detecting *out-of-distribution* samples under the open-world assumption. However, subtle differences exist among the sub-topics in terms of the specific definition and properties of ID and OOD data—which are often overlooked by the research community. To this end, we provide a clear introduction and description of each sub-topic in respective subsections (from Section 2.1 to 2.5). Each subsection details the motivation, background, formal definition, as well as relative position within the unified framework. Applications and benchmarks are also introduced, with concrete examples that facilitate understanding. Fig. 2 illustrates the settings for each sub-topic. In the end, we conclude this section by introducing the

neighborhood topics to clarify the scope of the generalized OOD detection framework. (Section 2.6).

Preliminary: Distribution Shift Key to our framework, the notion of distribution shift is very broad and can exhibit in various forms. There are two general types of distribution shift: covariate shift and semantic (label) shift. Formally, let \mathcal{X} be the input (sensory) space and \mathcal{Y} be the label (semantic) space, a data distribution is defined as a joint distribution $P(X, Y)$ over $\mathcal{X} \times \mathcal{Y}$. Distribution shift can occur in either the marginal distribution $P(X)$, or both $P(Y)$ and $P(X)$. Note that shift in $P(Y)$ naturally triggers shift in $P(X)$.

Examples of covariate distribution shift on $P(X)$ include adversarial examples [31], [32], domain shift [33], and style changes [34]. Importantly, we note that covariate shift is more commonly used to evaluate model *generalization* and robustness performance, where the label space \mathcal{Y} remains the same during test time. On the other hand, the detection of semantic distribution shift (*e.g.*, due to the occurrence of new classes) is the focal point of many *detection* tasks considered in this framework, where the label space \mathcal{Y} can be different between ID and OOD data and hence the model should not make any prediction.

With the concept of distribution shift in mind, readers can get a general idea of the differences and connections among sub-topics/tasks in Fig. 1. Notice that different sub-tasks can be easily identified with the following four dichotomies: 1) covariate / semantic shift dichotomy; 2) single / multiple class dichotomy; 3) ID classification needed / non-needed dichotomy; 4) inductive / transductive dichotomy. Next, we proceed with elaborating on each sub-topic.

2.1 Anomaly Detection

Background The notion of “anomaly” stands in contrast with the “normal” defined in advance. The concept of “normal” should be clear and reflect the real task. For example, to create a “not-hotdog detector”, the concept of the normal should be clearly defined as the hotdog class, *i.e.*, a food category, so that objects that violate this definition are identified as anomalies, which include steaks, rice, and non-food objects like cats and dogs. Ideally, “hotdog” would be regarded as a homogeneous concept, regardless of the sub-classes of French or American hotdog.

Current anomaly detection settings often restrict the environment of interest to some specific scenarios. For example, the “not-hotdog detector” only focuses on realistic images, assuming the nonexistence of images from other domains such as sketches. Another realistic example is industrial defect detection, which bases on only one set of assembly lines for a specific product. In other words, the “open-world” assumption is usually not completely “open”. Nevertheless, “not-hotdog” or “defects” can form a large unknown space that breaks the “closed-world” assumption.

In summary, the key to anomaly detection is to define normal clearly (usually without sub-classes) and detect all possible anomalous samples under some specific scenarios.

Definition Anomaly detection (AD) aims to detect any anomalous samples that are deviated from the predefined normality during testing. The deviation can happen due to either covariate shift or semantic shift, which leads to two sub-tasks: sensory AD and semantic AD, respectively [20].

Sensory AD detects test samples with covariate shift, under the assumption that normalities come from the same covariate distribution. No semantic shift takes place in sensory AD settings. On the other hand, semantic AD detects test samples with label shift, assuming that normalities come from the same semantic distribution (category), *i.e.*, normalities should belong to only one class.

Formally, in sensory AD, normalities are from in-distribution $P(X)$ while anomalies encountered at test time are from out-of-distribution $P'(X)$, where $P(X) \neq P'(X)$ — only covariate shift occurs. The goal in sensory AD is to detect samples from $P'(X)$. No semantic shift occurs in this setting, *i.e.*, $P(Y) = P'(Y)$. Conversely, for semantic AD, only semantic shift occurs (*i.e.*, $P(Y) \neq P'(Y)$) and the goal is to detect samples that belong to novel classes.

Remark: Sensory / Semantic Dichotomy Our sensory / semantic dichotomy for the AD sub-task definition comes from the low-level sensory anomalies and high-level semantic anomalies that are introduced in [35] and highlighted in the recent AD survey [20], to reflect the rise of deep learning. Note that although most sensory and semantic AD methods are shown to be mutually inclusive due to the common shift on $P(X)$, some approaches are specialized in one of the sub-tasks (ref. Section 5). Recent research communities are also trending on subdividing types of anomalies in order to develop targeted methods, so that practitioners can select the optimal solution for their own practical problem [35], [36].

Position in Framework Under the generalized OOD detection framework, the definition of “normality” seamlessly connects to the notion of “in-distribution”, and “anomaly” corresponds to “out-of-distribution”. Importantly, AD treats ID samples as a whole, which means that regardless of the number of classes (or statistical modalities) in ID data, AD does not require the differentiation in the ID samples. This feature is an important distinction between AD and other sub-topics such as OSR and OOD detection.

Application and Benchmark Sensory AD only focuses on objects with the same or similar semantics, and identifies the observational differences on their surface. Samples with sensory differences are recognized as sensory anomalies. Example applications include adversarial defense [37], forgery recognition of biometrics and artworks [38], [39], [40], [41], image forensics [42], [43], [44], industrial inspection [45], [46], [47], *etc.* The most popular academic AD benchmark is MVTec-AD [45] for industrial inspection.

In contrast to sensory AD, semantic AD only focuses on the semantic shift. An example of real-world applications is crime surveillance [48], [49]. Active image crawlers for a specific category also need semantic AD methods to ensure the purity of the collected images [50]. An example of the academic benchmarks is to recursively use one class from MNIST as ID during training, and ask the model to distinguish it from the rest of the 9 classes during testing.

Evaluation In the AD benchmarks, test samples are annotated to be either normal or abnormal. The deployed anomaly detector will produce a confidence score for a test sample, indicating how confident the model considers the sample as normality. Samples below the predefined confidence threshold are considered abnormal. By viewing the true normalities as positive and anomalies as negative, different

thresholds will produce a series of true positive rates (TPR) and false-positive rates (FPR)—from which we can calculate the area under the receiver operating characteristic curve (AUROC) [51]. Similarly, the precision and recall values can be used to compute metrics of F-scores and the area under the precision-recall curve (AUPR) [52]. Note that there can be two variants of AUPR values: one treating “normal” as the positive class, and the other treating “abnormal” as the positive class. For AUROC and AUPR, a higher value indicates better detection performance.

Remark: Alternative Taxonomy on Anomalies Some previous literature considers anomalies types to be three-fold: point anomalies, conditional or contextual anomalies, and group or collective anomalies [20], [21], [23]. In this survey, we mainly focus on point anomalies detection for its popularity in practical applications and its adequacy to elucidate the similarities and differences between sub-tasks. Details of the other two kinds of anomalies, *i.e.*, contextual anomalies that often occur in time-series tasks, and collective anomalies that are common in the data mining field, are not covered in this survey. We recommend readers to the recent AD survey papers [20] for in-depth discussion on them.

Remark: Taxonomy based on Supervision We use sensory / semantic dichotomy to subdivide AD at the task level. From the perspective of methodologies, some literature categorizes AD techniques into unsupervised and (semi-)supervised settings. Note that these two taxonomies are orthogonal as they focus on tasks and methods respectively.

2.2 Novelty Detection

Background The word “novel” generally refers to the unknown, new, and something interesting. While novelty detection (ND) is often interchangeable with AD in the community, strictly speaking, their subtle difference is worth noticing. In terms of motivation, novelty detection usually does not perceive “novel” test samples as erroneous, fraudulent, or malicious as AD does, but cherishes them as learning resources for potential future use with a positive learning attitude [20], [21]. In fact, novelty detection is also known as “novel class detection” [26], [27], indicating that it is primarily focusing on detecting semantic shift.

Definition Novelty detection aims to detect any test samples that do not fall into any training category. The detected novel samples are usually prepared for future constructive procedures, such as more specialized analysis, or incremental learning of the model itself. Based on the number of training classes, ND contains two different settings: 1) one-class novelty detection (*one-class ND*): only one class exists in the training set; 2) multi-class novelty detection (*multi-class ND*): multiple classes exist in the training set. It is worth noting that despite having many ID classes, the goal of multi-class ND is only to distinguish novel samples from ID. Both one-class and multi-class ND are formulated as binary classification problems.

Position in Framework Under the generalized OOD detection framework, ND deals with the setting where OOD samples have semantic shift, without the need for classification in the ID set even if possible. Therefore, ND shares the same problem definition with semantic AD.

Application and Benchmark Real-world ND application includes video surveillance [48], [49], planetary exploration [53] and incremental learning [54], [55]. For one-class ND, an example academic benchmark can be identical to that of semantic AD, which considers one class from MNIST as ID and the rest as the novel. The corresponding MNIST benchmark for multi-class ND may use the first 6 classes during training, and test on the remaining 4 classes as OOD.

Evaluation The evaluation of ND is identical to AD, which is based on AUROC, AUPR, or F-scores (see details in Section 2.1).

Remark: One-Class / Multi-Class Dichotomy Although the ND models do not require the ID classification even with multi-class annotations, the method on multi-class ND can be different from one-class ND, as multi-class ND can make use of the multi-class classifier while one-class ND cannot. Also note that semantic AD can be further split into one-class semantic AD and multi-class semantic AD that matches ND, as semantic AD is equivalent to ND.

Remark: Nuance between AD and ND Apart from the special interest in semantics, some literature [56], [57] also point out that ND is supposed to be fully unsupervised, while AD might have some abnormal training samples.

2.3 Open Set Recognition

Background Machine learning models trained in the closed-world setting can incorrectly classify test samples from unknown classes as one of the known categories with high confidence [58]. Some literature refers to this notorious overconfident behavior of the model as “arrogance”, or “agnostophobia” [59]. Open set recognition (OSR) is proposed to address this problem, with their own terminology of “known known classes” to represent the categories that exist at training, and “unknown unknown classes” for test categories that do not fall into any training category. Some other terms, such as open category detection [60] and open set learning [61], are simply different expressions for OSR.

Definition Open set recognition requires the multi-class classifier to simultaneously: 1) accurately classify test samples from “known known classes”, and 2) detect test samples from “unknown unknown classes”.

Position in Framework OSR well aligns with our generalized OOD detection framework, where “known known classes” and “unknown unknown classes” correspond to ID and OOD respectively. Formally, OSR deals with the case where OOD samples during testing have semantic shift, *i.e.*, $P(Y) \neq P'(Y)$. The goal of OSR is largely shared with that of multi-class ND—the only difference is that OSR additionally requires accurate classification of ID samples from $P(Y)$.

Application and Benchmark OSR supports the robust deployment of real-world image classifiers in general, which can reject unknown samples in the open world [62], [63]. An example academic benchmark on MNIST can be identical to multi-class ND, which considers the first 6 classes as ID and the remaining 4 classes as OOD. In addition, OSR further requires a good classifier on the 6 ID classes.

Evaluation Similar to AD and ND, the metrics for OSR include F-scores, AUROC, and AUPR. Beyond them, the classification performance is also evaluated by standard ID

accuracy. While the above metrics evaluate the novelty detection and ID classification capabilities independently, some works raise some evaluation criteria for joint evaluation, such as CCR@FPR x [59], which calculates the class-wise recall when a certain FPR equal to x (e.g., 10^{-1}) is achieved.

2.4 Out-of-Distribution Detection

Background With the observation that deep learning models are often inappropriate but in fact overconfident in classifying samples from different semantic distributions in the image classification task and text categorization [64], the field of out-of-distribution detection emerges, requiring the model to reject inputs that are semantically different from the training distribution and therefore should not be predicted by the model.

Definition Out-of-distribution detection, or OOD detection, aims to detect test samples that drawn from a distribution that is different from the training distribution, with the definition of distribution to be well-defined according to the application in the target. For most machine learning tasks, the distribution should refer to “label distribution”, which means that OOD samples should not have overlapping labels w.r.t. training data. Formally, in the OOD detection, the test samples come from a distribution whose semantics are shifted from ID, i.e., $P(Y) \neq P'(Y)$. Note that the training set usually contains multiple classes, and OOD detection should NOT harm the ID classification capability.

Position in Framework Out-of-distribution detection can be canonical to OSR in common machine learning tasks like multi-class classification—keeping the classification performance on test samples from ID class space \mathcal{Y} , and reject OOD test samples with semantics outside the support of \mathcal{Y} . Also, the multi-class setting and the requirement of ID classification distinguish the task from AD and ND.

Application and Benchmark The application of OOD detection usually falls into safety-critical situations, such as autonomous driving [65], [66]. An example academic benchmark is to use CIFAR-10 as ID during training, and to distinguish CIFAR images from other datasets such as SVHN, etc. Researchers should pay attention that OOD datasets should NOT have label overlapping with ID datasets when building the benchmark.

Evaluation Apart from F-scores, AUROC, and AUPR, another commonly-used metric is FPR@TPRx, which measures the FPR when the TPR is x (e.g., 0.95). Some works also use an alternative metric, TNR@TPRx, which is equivalent to 1-FPR@TPRx. OOD detection also concerns the performance of ID classification.

Remark: OSR vs. OOD Detection The difference between OSR and OOD detection tasks are three-fold.

1) Different benchmark setup: OSR benchmarks usually split one multi-class classification dataset into ID and OOD parts according to classes, while OOD detection takes one dataset as ID and find several other datasets as OOD with the guarantee of non-overlapping categories between ID / OOD datasets. However, despite the different benchmark traditions of the two sub-tasks, they are in fact tackling the same problem of semantic shift detection.

2) No additional data in OSR: Due to the requirement of theoretical open-risk bound guarantee, OSR discourages the usage of additional data during training by design [28]. This restriction precludes methods that are more focused on effective performance improvements (e.g., outlier exposures [67]) but may violate OSR constraints.

3) Broadness of OOD detection: Compare to OSR, OOD detection encompasses a broader spectrum of learning tasks (e.g., multi-label classification [68]), wider solution space (to be discussed in Section 3).

Remark: Mainstream OOD Detection Focuses on Semantics While most works in the current community interpret the keyword “out-of-distribution” as “out-of-label/semantic-distribution”, some OOD detection works also consider detecting covariate shifts [69], which claim that covariate shift usually leads to a significant drop in model performance and therefore needs to be identified and rejected. However, although detecting covariate shift is reasonable on some specific tasks (usually due to high-risk or privacy reasons) that to be discussed in the following paragraph, research on this topic remains a controversial task *w.r.t* OOD generalization tasks (*c.f.* Section 2.6 and Section 7.2). Detecting semantic shift has been the mainstream of OOD detection tasks.

Remark: Extension of OOD Detection We provide another definition from the perspective of generalization: Out-of-distribution detection, or OOD detection, aims to detect test samples to which the model cannot or does not want to generalize [70]. In the most of the machine learning tasks, such as image classification, the models are expected to generalize their prediction capability to samples with covariate shift, and they are only unable to generalize when semantic shift occurs. However, for applications where models are by-design nontransferable to other domain, such as many deep reinforcement learning tasks like game AI [71], [72], the key term “distribution” should refer to “data / input distribution”, so that the model should refuse to make decision under the environment that is not the same as training environment, i.e., $P(X) \neq P'(X)$. Similar applications are those high-risk tasks such as medical image classification [73] or in privacy-sensitive scenario [74], where the models are expected to be very conservative and only make predictions for samples exactly from the training distribution, rejecting any samples that deviate from it. In sum, an OOD detection task can be valid when the “detection” task has well reconciled with the “generalization” task. By all means, detecting semantic shift is still the mainstream of OOD detection tasks and is the focus of this survey.

2.5 Outlier Detection

Background According to Wikipedia [75], an outlier is a data point that differs significantly from other observations. Recall that the problem settings in AD, ND, OSR, and OOD detect unseen test samples that are different from the training data distribution. In contrast, outlier detection directly processes all observations and aims to select outliers from the contaminated dataset [16], [17], [18]. Since outlier detection does not follow the train-test procedure but has access to all observations, approaches to this problem are usually transductive rather than inductive [76].

Definition Outlier detection aims to detect samples that are markedly different from the others in the given observation set, due to either covariate or semantic shift.

Position in Framework Different from all previous sub-tasks, whose in-distribution is defined during training, the “in-distribution” for outlier detection refers to the majority of the observations. Outliers may exist due to semantic shift on $P(Y)$, or covariate shift on $P(X)$.

Application and Benchmark While mostly applied in data mining tasks [77], [78], [79], outlier detection is also used in the real-world computer vision applications such as video surveillance [80] and dataset cleaning [81], [82], [83]. For the application of dataset cleaning, outlier detection is usually used as a pre-processing step for the main tasks such as learning from open-set noisy labels [84], webly supervised learning [85], and open-set semi-supervised learning [86]. To construct an outlier detection benchmark on MNIST, one class should be chosen so that all samples that belong to this class are considered as inliers. A small fraction of samples from other classes are introduced as outliers to be detected.

Evaluation Apart from F-scores, AUROC, and AUPR, the evaluation of outlier detectors can be also evaluated by the performance of the main task it supports. For example, if an outlier detector is used to purify a dataset with noisy labels, the performance of a classifier that is trained on the cleaned dataset can indicate the quality of the outlier detector.

Remark: On Inclusion of Outlier Detection Interestingly, the outlier detection task can be considered as an outlier in the generalized OOD detection framework, since outlier detectors are operated on the scenario when all observations are given, rather than following the training-test scheme. Also, publications exactly on this topic are rarely seen in the recent deep learning venues. However, we still include outlier detection in our framework, because intuitively speaking, outliers also belong to one type of out-of-distribution, and introducing it can help familiarize readers more with various terms (*e.g.*, OD, AD, ND, OOD) that have confused the community for a long while.

2.6 Related Topics

Apart from the five sub-topics that are described in our *generalized OOD detection* framework (shown in Figure 1), we further briefly discuss five related topics below, which helps clarify the scope of this survey.

Learning with Rejection (LWR) LWR can date back to early works on abstention [87], [88], which considered simple model families such as SVMs [89]. The phenomenon of neural networks’ overconfidence in OOD data is first revealed by [90]. Despite methodologies differences, subsequent works developed on OOD detection and OSR share the underlying spirit of classification with rejection option.

Domain Adaptation / Generalization Domain Adaptation (DA) [15] and Domain Generalization (DG) [91] also follow “open-world” assumption. Different from generalized OOD detection settings, DA/DG expects the existence of covariate shift during testing without any semantic shift, and requires classifiers to make accurate predictions into the same set of classes [92]. Noticing that OOD detection commonly concerns detecting the semantic shift, which is

complementary to DA/DG. In the case when both covariate and semantic shift take place, the model should be able to detect semantic shift while being robust to covariate shift. More discussion on relations between DA/DG and OOD detection is in Section 7.2. The difference between DA and DG is that while the former requires extra but few training samples from the target domain, the latter one does not.

Novelty Discovery Novelty discovery [93], [94], [95], [96], [97] requires all observations are given in advance as outlier detection does. The observations are provided in a *semi-supervised manner*, and the goal is to explore and discover the new categories and classes in the unlabeled set. Different from outlier detection where outliers are sparse, the unlabeled set in novelty discovery setting can mostly consist of, and even be overwhelmed by unknown classes.

Zero-shot Learning Zero-shot learning [98] has a similar goal of novelty discovery, but follows the *training-testing scheme*. The test set is under the “open-world” assumption with unknown classes, which expect classifiers trained only on the known classes to perform classification on unknown testing samples with the help of extra information such as label relationships.

Open-world Recognition Open-world recognition [99] aims to build a lifelong learning machine that can actively detect novel images [100], label them as new classes, and perform continuous learning. It can be viewed as the combination of novelty detection and incremental learning.

2.7 Organization of Remaining Sections

In the following sections, we will introduce the methodologies of AD, ND, OSR, OOD detection, and OD, each sub-task in one section. For each sub-task, we categorize and introduce the methodologies into the following four groups: 1) **classification-based methods**: methods that largely rely on classifiers; 2) **density-based methods**: detecting OOD by modeling data density; 3) **distance-based methods**: using distance metrics (usually in the *feature space*) to identify OODs; and 4) **reconstruction-based methods**: methods featured by reconstruction techniques. Our following review mainly focuses on summarizing the main ideas of the methodological categories, hoping to provide readers with big pictures for solving OOD-related problems and inspire readers to develop more effective methods. A concurrent survey [101] on reviewing methods under the generalized OOD detection framework provides a detailed explanation of OOD-related methods, greatly complementing our work.

Remark: This Survey Highlights OOD Detection Methods Notice the existence of comprehensive surveys on AD, ND, and OD methodologies in recent years [20], [21], [22], [23], this survey will highlight the overview of OSR and OOD detection methods in complementing the good surveys mentioned above, so that the ordering of the following methodological sections are: OOD detection in Section 3, OSR in Section 4, AD / ND in Section 5, and OD in Section 6.

3 OOD DETECTION: METHODOLOGY

In this section, we introduce the methodology for OOD detection. We first present *classification-based* model in Section 3.1, followed by density-based methods in Section 3.2.

TABLE 1: Paper list for generalized out-of-distribution detection tasks.

Task		Methodology		Reference
§ 3 Out-of-Distribution Detection	Classification	§ 3.1 Output-based Methods	a: Post-hoc Detection	[64], [102], [103], [104], [105], [106], [107], [108], [109], [110], [111]
			b: Conf. Enhancement	[69], [112], [113], [114], [115], [116], [117], [118], [119], [120], [121], [122], [123], [124], [125]
			c: Outlier Exposure	[59], [67], [126], [127], [128], [129], [130], [131], [132], [133], [134], [135], [136]
		§ 3.1.2: Label Space Redesign		[137], [138], [139], [140], [141]
		§ 3.1.3: OOD Data Generation		[142], [143], [144], [145], [146], [147]
		§ 3.1.4: Gradient-based Methods		[102], [148]
	§ 3.2: Density-based Methods	§ 3.1.5: Bayesian Models		[149], [150], [151], [152], [153], [154]
		§ 3.1.6: Large-Scale OOD Detection		[138], [140], [155], [156]
		§ 3.2: Density-based Methods		[103], [157], [158], [159], [160], [161], [162], [163], [164], [165], [166], [167], [168], [169], [170], [171], [172]
		§ 3.3: Distance-based Methods		[103], [110], [173], [174], [175], [176], [177], [178], [179]
	§ 3.4: Reconstruction-based Methods			[180], [181], [182]
§ 4 Open Set Recognition	Classification	§ 4.1.1: EVT-based Calibration		[183], [184], [185], [186]
		§ 4.1.2: EVT-free Calibration		[187], [188], [189]
		§ 4.1.3: Unknown Generation		[190], [191], [192], [193], [194], [195], [196]
	Reconstruction	§ 4.2: Distance-based Methods		[197], [198], [199], [200], [201], [202], [203]
		§ 4.3.1: Sparse Representation		[204], [205], [206]
		§ 4.3.2: Reconstruction-Error		[207], [208], [209], [210]
§ 5 Anomaly Detection & Novelty Detection	Density	§ 5.1.1: Classic Density Est.		[211], [212], [213], [214], [215], [216]
		§ 5.1.2: NN-based Density Est.		[157], [158], [159], [159], [160], [162], [217], [218], [219], [220], [221], [222], [223]
		§ 5.1.3: Energy-based Models		[224], [225], [226]
		§ 5.1.4: Frequency-based Methods		[227], [228], [229], [230]
	Reconstruction	§ 5.2.1: Sparse Representation		[231], [232], [233], [234], [235]
		§ 5.2.2: Reconstruction-Error		[56], [159], [161], [236], [236], [237], [237], [238], [238], [239], [240], [241], [242], [243], [244], [245], [246], [247]
		§ 5.2.3: Gradient Representation		[248]
	Classification	§ 5.3: Distance-based Methods		[249], [250], [251]
		§ 5.4.1: One-Class Classification		[252], [253], [254], [255]
		§ 5.4.2: PU Learning		[256], [257], [258], [259], [260], [261], [262], [263], [264], [265], [266], [267], [268]
		§ 5.4.3: Self-Supervised Learning		[269], [270], [271], [272], [273]
	§ 5.5: Discussion and Theoretical Analysis			[60], [61]
§ 6 Outlier Detection	Distance	§ 6.1: Density-based Methods		[274], [275], [276], [277], [278], [279], [280], [281]
		§ 6.2.1: Cluster-based Methods		[282], [283]
		§ 6.2.2: Graph-based Methods		[284], [285], [286], [287], [288]
	§ 6.3: Classification-based Methods			[252], [253], [269], [289], [290], [291], [292]

Distance-based methods will be introduced in Sections 3.3. A brief discussion will be included at the end.

3.1 Classification-based Methods

Research on OOD detection originated from a simple baseline, that is, using the maximum softmax probability as the indicator score of ID-ness [64]. Early OOD detection methods focus on deriving improved OOD scores based on the output of neural networks.

3.1.1 Output-based Methods

a. Post-hoc Detection Post-hoc methods have the advantage of being easy to use without modifying the training procedure and objective. The property can be important for the adoption of OOD detection methods in real-world production environments, where the overhead cost of re-training can be prohibitive. Early work ODIN [102] is a post-hoc method that uses temperature scaling and input perturbation to amplify the ID/OOD separability. Key to the method, a sufficiently large temperature has a strong

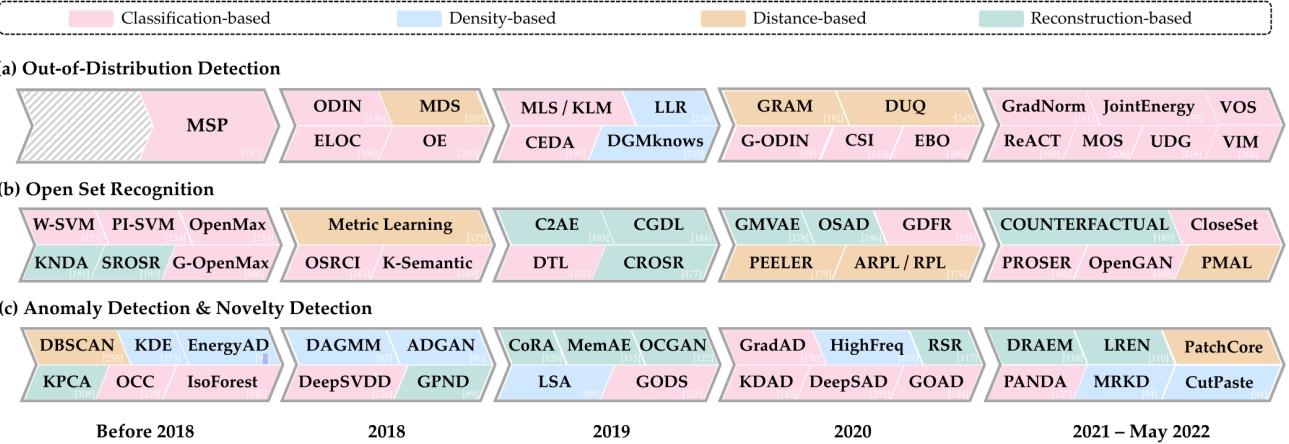


Fig. 3: Timeline for representative methodologies of (a) OOD detection, details in Section 3, (b) open set recognition, details in Section 4, and (c) anomaly detection & novelty detection, details in Section 5, under generalized OOD detection framework. Different colors indicate different categories of methodologies. Each method has its corresponding reference (inconspicuous white) in the lower right corner. Methods with open-source code are prioritized for inclusion in this figure.

smoothing effect that transforms the softmax score back to the logit space—which effectively distinguishes ID vs. OOD. Note that this is different from confidence calibration, where a much milder T is employed. While calibration focuses on representing the true correctness likelihood of ID data only, the ODIN score is designed to maximize the gap between ID and OOD data and may no longer be meaningful from a predictive confidence standpoint. Built on the insights, recent work [104], [111] proposed using an energy score for OOD detection, which enjoys theoretical interpretation from a likelihood perspective [293]. Test samples with lower energy are considered ID and vice versa. JointEnergy score [106] is then proposed to perform OOD detection for multi-label classification networks. Techniques such as layer-wise Mahalanobis distance [103] and Gram Matrix [105] are implemented for better hidden feature quality to perform density estimation.

Recently, one fundamental cause of the overconfidence issue on OOD data is revealed that using mismatched BatchNorm statistics—that are estimated on ID data yet blindly applied to the OOD data in testing—can trigger abnormally high unit activations and model output accordingly [107]. Therefore, ReAct [107] proposes truncating the high activations, which establishes strong post-hoc detection performance and further boosts the performance of existing scoring functions. Similarly, NMD [108] uses the activation means from BatchNorm layers for ID / OOD discrepancy. While ReAct considers activation space, [109] proposes a weight sparsification-based OOD detection framework termed DICE. DICE ranks weights based on a measure of contribution, and selectively uses the most salient weights to derive the output for OOD detection. By pruning away noisy signals, DICE provably reduces the output variance for OOD data, resulting in a sharper output distribution and stronger separability from ID data.

b. Confidence Enhancement Methods Tailored for OOD detection problem, confidence can be developed via designing a confidence-estimating branch [112] or class [113],

data augmentation [114], ensembling with leaving-out strategy [115], adversarial training [116], [117], [118], [119], [129], stronger data augmentation [120], [121], [122], [123], better uncertainty modeling [124], [294], and utilizing feature or statistics from the intermediate-layer features [111], [295]. Specially, to enhance the sensitivity to covariate shift, some methods focus on the hidden representations in the middle layers of neural networks. Generalized ODIN, or G-ODIN [69] extended ODIN [102] by using a specialized training objective termed DeConf-C and choose hyperparameters such as perturbation magnitude on ID data. Note that we do not categorize G-ODIN as post-hoc method as it requires model retraining. Recent work [125] shows that the overconfidence issue can be mitigated through Logit Normalization (LogitNorm), a simple fix to the common cross-entropy loss by enforcing a constant vector norm on the logits in training. Trained with LogitNorm, neural networks produce highly distinguishable confidence scores between in- and out-of-distribution data.

c. Outlier Exposure Another branch of OOD detection methods makes use of a set of collected OOD samples, or “outlier”, during training to help models learn ID/OOD discrepancy. Starting from the concurrent baselines that encourage a flat/high-entropic prediction on given OOD samples [59], [67] and suppressing OOD feature magnitudes [59], a follow-up work, MCD [126] uses a network with two branches, between which entropy discrepancy is enlarged for OOD training data. Another straightforward approach with outlier exposure spares an extra abstention (or rejection class) and considers all the given OOD samples in this class [129], [132], [134]. A later work OECC [127] noticed that an extra regularization for confidence calibration introduces additional improvement for OE. To effectively utilize the given, usually massive, OOD samples, some works use outlier mining [129], [136] and adversarial resampling [128] approaches to obtain a compact yet representative set. Other works consider a more practical scenario where given OOD samples contain ID samples, therefore using pseudo-labeling [130]

or ID filtering methods [131] to reduce the interference of ID data. In general, OOD detection with outlier exposure can reach a much better performance. However, as research shows that the performance can be largely affected by the correlations between given and real OOD samples [133]. To address the issue, recent work [135] proposes a novel framework that enables effectively exploiting unlabeled in-the-wild data for OOD detection. Unlabeled wild data is frequently available since it is produced essentially for free whenever deploying an existing classifier in a real-world system. This setting can be viewed as training OOD detectors in their *natural habitats*, which provide a much better match to the true test time distribution than data collected offline.

3.1.2 Label Space Redesign

One-hot encoding is commonly used to encode categorical information for classification. However, one-hot encoding ignores the inherent relationship among labels. For example, it is **unreasonable** to have a **uniform distance** between **dog** and **cat** vs. **dog** and **car**. To this end, several works attempt to use information in the label space for OOD detection. Some works arrange the large semantic space into a hierarchical taxonomy of known classes [137], [138]. Under the redesigned label architecture, top-down classification strategy [137] and group softmax training [138] are demonstrated effective. Another set of works uses word embeddings to automatically construct the label space. In [139], the sparse one-hot labels are replaced with several dense word embeddings from different NLP models, forming multiple regression heads for robust training. When testing, the label, which has the minimal distance to all the embedding vectors from different heads, will be considered as the prediction. If the minimal distance crosses above the threshold, the sample would be classified as “novel”. Recent works further take the image features from language-image pre-training models [296] to better detect novel classes, where the image encoding space also contains rich information from the language space [140], [141].

3.1.3 OOD Data Generation

The outlier exposure approaches impose a strong assumption on the availability of OOD training data, which can be infeasible in practice. When no OOD sample is available, some methods attempt to **synthesize OOD samples to enable ID/OOD separability**. Existing works leverage GANs to generate OOD training samples and force the model predictions to be uniform [142], generate boundary samples in the low-density region [143], or similarly, high-confidence OOD samples [144]. However, synthesizing images in the high-dimensional pixel space can be difficult to optimize. Recent work VOS [146] proposed synthesizing virtual outliers from the low-likelihood region in the feature space, which is more tractable given lower dimensionality. In object detection, [147] proposes synthesizing unknown objects from videos in the wild using spatial-temporal unknown distillation.

3.1.4 Gradient-based Methods

Existing OOD detection approaches primarily rely on the **output** (Section 3.1) or **feature space** for deriving OOD scores, while overlooking information from the **gradient**

space. ODIN [102] first explored using gradient information for OOD detection. In particular, ODIN proposed using input pre-processing by adding small perturbations obtained from the input gradients. The goal of ODIN perturbations is to increase the softmax score of any given input by reinforcing the model’s belief in the predicted label. Ultimately the perturbations have been found to create a greater gap between the softmax scores of ID and OOD inputs, thus making them more separable and improving the performance of OOD detection. While ODIN only uses gradients implicitly through input perturbation, recent work proposed GradNorm [148] which explicitly derives a scoring function from the gradient space. GradNorm employs the vector norm of gradients, backpropagated from the KL divergence between the softmax output and a uniform probability distribution.

3.1.5 Bayesian Models

A Bayesian model is a **statistical model** that implements Bayes’ rule to infer all **uncertainty** within the model [297]. The most representative method is the Bayesian neural network [298], which draws samples from the posterior distribution of the model via MCMC [299], Laplace methods [300], [301] and variational inference [302], forming the epistemic uncertainty of the model prediction. However, their obvious shortcomings of inaccurate predictions [303] and high computational costs [304] prevent them from wide adoption in practice. Recent works attempt several less principled approximations including MC-dropout [149] and deep ensembles [150], [305], [306] for faster and better estimates of uncertainty. These methods are less competitive for OOD uncertainty estimation. Further exploration takes natural-gradient variational inference and enables practical and affordable modern deep learning training while preserving the benefits of Bayesian principles [151]. Dirichlet Prior Network (DPN) is also used for OOD detection with an uncertainty modeling of three different sources of uncertainty: model uncertainty, data uncertainty, and distributional uncertainty and form a line of works [152], [153], [154]. Recently, Bayesian hypothesis test is used to formulate OOD detection, with upweighting method and Hessian approximation for scalability [307].

3.1.6 Large-scale OOD Detection

Recent works have advocated for OOD detection in large-scale settings, which are closer to real-world applications. Research efforts include scaling OOD detection to **large semantic label space** and exploiting **large pre-trained models**. For example, [138] revealed that approaches developed on the CIFAR benchmark might not translate effectively into ImageNet benchmark with a large semantic space, highlighting the need to evaluate OOD detection in a large-scale real-world setting. To overcome the challenge, the key idea of MOS [138] is to decompose the large semantic space into smaller groups with similar concepts, which allows simplifying the decision boundaries between known vs. unknown data. Recently, powerful pre-trained models have achieved astonishing results on various tasks and modalities. Several concurrent works [140], [155], [156] demonstrate that strong pretrained transformers can significantly improve some particularly difficult OOD tasks.

3.2 Density-based Methods

Density-based methods in OOD detection explicitly model the in-distribution with some probabilistic models, and flag test data in low-density regions as OOD. Although OOD detection can be different from AD in that multiple classes exist in the in-distribution, density estimation methods used for AD in Section 5.1.2 can be directly adapted to OOD detection by unifying the ID data as a whole [157], [158], [159], [160], [161]. When the ID contains multiple classes, class-conditional Gaussian distribution can explicitly model the in-distribution so that the OOD samples can be identified based on their likelihoods [103]. Flow-based methods [162], [163], [164], [165], [166] can also be used for probabilistic modeling. While directly estimating the likelihood seems like a natural approach, some works [167], [168], [169] find that probabilistic models sometimes assign a higher likelihood for the OOD sample. Several works attempt to solve the problems using likelihood ratio [170]. [171] finds that the likelihood exhibits a strong bias towards the input complexity and proposes a likelihood ratio-based method to compensate the influence of input complexity. Recent methods turn to new scores such as likelihood regret [172] or an ensemble of multiple density models [168]. To directly model the density of semantic space, SEM score is used with a simple combination of density estimation in the low-level and high-level space [308]. Overall, generative models can be prohibitively challenging to train and optimize, and the performance can often lag behind the classification-based approaches (Section 3.1).

3.3 Distance-based Methods

The basic idea of distance-based methods is that the testing OOD samples should be relatively far away from the centroids or prototypes of in-distribution classes. [103] uses the minimum Mahalanobis distance to all class centroids for detection. A subsequent work splits the images into foreground and background, and then calculates the Mahalanobis distance ratio between the two spaces [173]. In contrast to parametric approach, recent work [110] shows strong promise of non-parametric nearest-neighbor distance for OOD detection. Unlike Mahalanobis, the non-parametric approach does not impose any distributional assumption about the underlying feature space, hence providing stronger simplicity, flexibility and generality.

Some works use cosine similarity between test sample features and class features to determine OOD samples [174], [175]. The one-dimensional subspace spanned by the first singular vector of the training features is shown to be more suitable for cosine similarity-based detection [176]. Moreover, other works leverage distances with radial basis function kernel [177], Euclidean distance [178], and geodesic distance [309] between the input's embedding and the class centroids. Apart from calculating the distance between samples and class centroids, feature norm that in the orthogonal complement space of the principal space is shown effective on OOD detection [310]. Distance-based OOD scores can generally benefit from enhanced representations. Recent work CIDER [179] introduces a new representation learning framework for OOD detection, which consists of a dispersion loss promoting large angular distances among different

class prototypes, along with a compactness loss encouraging samples to be close to their class prototypes.

3.4 Reconstruction-based Methods

The core idea of reconstruction-based methods is that the encoder-decoder framework trained on the ID data usually yields different outcomes for ID and OOD samples. The difference in model performance can be utilized as an indicator for detecting anomalies. For example, reconstruction models that are only trained by ID data cannot well recover the OOD data [180], and therefore the OOD can be identified. While reconstruction-based models with pixel-level comparison seem not a popular solution in OOD detection for its expensive training cost, reconstructing with hidden features is shown as a promising alternative [181]. Rather than reconstructing the entire image, recent work MoodCat [182] masks a random portion of the input image and identifies OOD samples using the quality of the classification-based reconstruction results.

3.5 Discussion

The field of OOD detection has enjoyed rapid development since its emergence, with a large space of solutions. In the multi-class setting, the problem can be canonical to OSR (Section 4)—accurately classify test samples from ID within the class space \mathcal{Y} , and reject test samples with semantics outside the support of \mathcal{Y} . The difference often lies in the evaluation protocol. OSR splits a dataset into two halves: one set as ID and another set as OOD. In contrast, OOD allows a more general and flexible evaluation by considering test samples from different datasets or domains. Moreover, OOD detection encompasses a broader spectrum of learning tasks (e.g., multi-label classification [106], object detection [146], [147]) and solution space. Apart from the methodology development, theoretical understanding has also received attention in the community [293], providing provable guarantees and empirical analysis to understand how OOD detection performance changes with respect to data distributions.

4 OPEN SET RECOGNITION: METHODOLOGY

In this section, we introduce the methodology for multi-class ND and open-set recognition (OSR) together. We discuss these two sub-tasks together since both tasks focus on the scenario where ID data comprises multiple classes. The only difference is that OSR has an additional objective to accurately classify the ID data, while multi-class ND produces an ID/OOD binary classifier.

Since multi-class ND and OSR consider multiple classes during training, most methods are classification-based (Section 4.1). Alternative methods are based on ID prototypes (Section 4.2) and reconstruction (Section 4.3). Few density-based methods will be discussed in Section 4.4 along with a discussion.

4.1 Classification-based Methods

The concept of OSR was first introduced in [58], which showed the validity of 1-class SVM and binary SVM for

solving the OSR problem. In particular, [58] proposes the 1-vs-Set SVM to manage the open-set risk by solving a two-plane optimization problem instead of the classic half-space of a binary linear classifier. This paper highlighted that the open-set space should also be bounded, in addition to bounding the ID risk.

4.1.1 EVT-based Uncertainty Calibration

Early works observe the overconfidence of neural networks and therefore focus on redistributing the logits by using the compact abating probability (CAP) [183] and extreme value theory (EVT) [184], [311], [312]. In particular, classic probabilistic models lack the consideration of open-set space. CAP explicitly models the probability of class membership abating from ID points to OOD points, and EVT focuses on modeling the tail distribution with extreme high/low values. In the context of deep learning, OpenMax [185] first implements EVT for neural networks. OpenMax replaces the softmax layer with an OpenMax layer, which calibrates the logits with a per-class EVT probabilistic model such as Weibull distribution. OpenMax also provides alternative solutions by using penultimate features for EVT modeling, forming a density-based method.

4.1.2 EVT-Free Confidence Enhancement

To circumvent the requirement of constructing open-set risks, some works achieved good empirical results without EVT. For example, [187] uses a membership loss to encourage high activations for known classes, and uses large-scale external datasets to learn globally negative filters that can reduce the activations of novel images. Apart from explicitly forcing discrepancy between known/unknown classes, other methods extract stronger features through an auxiliary task of transformation classification [188], or mutual information maximization between the input image and its latent features [189], etc.

4.1.3 Unknown Class Generation

Image generation techniques have been utilized to synthesize unknown samples from known classes, which helps distinguish between known vs. unknown samples [190], [191], [192], [193]. While these methods are promising on simple images such as handwritten characters, they do not scale to complex natural image datasets due to the difficulty in generating high-quality images in high-dimensional space. Another solution is to successively choose random categories in the training set and treat them as unknown, which helps the classifier to shrink the boundaries and gain the ability to identify unknown classes [194], [195]. Moreover, [196] splits the training data into typical and atypical subsets, which also helps learn compact classification boundaries.

4.2 Distance-based Methods

Distance-based methods for OSR require the prototypes to be class-conditional, which allows maintaining the ID classification performance. Category-based clustering and prototyping are performed based on the visual features extracted from the classifiers. OOD samples can be detected by computing the distance *w.r.t.* clusters [197], [198]. Some methods also leveraged contrastive learning to learn more

compact clusters for known classes [199], [200], which enlarge the distance between ID and OOD. CROSR [201] enhances the features by concatenating visual embeddings from both the classifier and reconstruction model for distance computation in the extended feature space. Besides using features from classifiers, GMVAE [202] extracts features using a reconstruction VAE, and models the embeddings of the training set as a Gaussian mixture with multiple centroids for the following distance-based operations. Classifiers using nearest neighbors are also adapted for OSR problem [203]. By storing the training samples, the nearest neighbor distance ratio is used for identifying unknown samples in testing.

4.3 Reconstruction-based Methods

With similar motivations as Section 3.4, reconstruction-based methods expect different reconstruction behavior for ID vs. OOD samples. The difference can be captured in the latent feature space or the pixel space of reconstructed images.

4.3.1 Sparse Representation Methods

By sparsely encoding images from the known classes, open-set samples can be identified based on their dense representation. Techniques such as sparsity concentration index [204] and kernel null space methods [205], [206] are used for sparse encoding.

4.3.2 Reconstruction-Error Methods

By fixing the visual encoder obtained from standard multi-class training to maintain ID classification performance, C2AE trains a decoder conditioned on label vectors and estimates the reconstructed images using EVT to distinguish unknown classes [207]. Subsequent works use conditional Gaussian distributions by forcing different latent features to approximate class-wise Gaussian models, which enables classifying known samples as well as reject unknown samples [208]. Other methods generate counterfactual images, which help the model focus more on semantics [209]. Adversarial defense is also considered in [210] to enhance model robustness.

4.4 Discussion

Although there is not an independent section for density-based methods, these methods can play an important role and are fused as a critical step in some classification-based methods such as OpenMax [185]. The density estimation on visual embeddings can effectively detect unknown classes without influencing the classification performance. A hybrid model also uses a flow-based density estimator to detect unknown samples [313].

Due to the restriction on using only ID data for training, OSR methods do not implement background classes, or outlier exposure (*c.f.* Section 3.1.1). A recent research shows that a good classifier [314] in the close-set is comparable to the state-of-the-art OSR methods, which urges the community to develop more general approaches without great efforts on carving for specific datasets.

5 ANOMALY DETECTION & NOVELTY DETECTION: METHODOLOGY

In this section, we review methodologies for AD and one-class ND. Most of the methods for sensory AD and semantic AD are shared, except for sensory AD focuses more on local information in the images and internal information of neural networks. Moreover, semantic AD and one-class ND have the same problem formulation (recall Section 2.2), therefore we review the methods for these three sub-tasks altogether.

Given the homogeneous in-distribution data, a straightforward approach is to estimate the in-distribution density and reject OOD test samples that deviate from the estimated distribution. We summarize density-based methods in Section 5.1. Alternative methods rely on the quality of image reconstruction to distinguish anomalous samples (Section 5.2), or directly learn a decision boundary between ID and OOD data (Section 5.4). We also review distance-based in Section 5.3. Lastly, we conclude with a discussion and present theoretical works in Section 5.5.

5.1 Density-based Methods

Density-based methods attempt to model the distribution of normal data (ID), with an operating assumption that anomalous test data has low likelihood whereas normal data has higher likelihood under the estimated density model.

5.1.1 Classic Density Estimation

a. Parametric Density Estimation Parametric density estimation assumes the ID density can be expressed through some pre-defined distributions [211]. One approach is to fits a multivariate Gaussian distribution on the training data and measures the Mahalanobis distance between the test sample and the expectation of training samples [212], [315]. Other works adopt more complex assumptions on in-distribution, such as mixed Gaussian distribution [213], [316], and Poisson distribution [214], etc.

b. Non-parametric Density Estimation Nonparametric density estimation solves a more practical scenario where a predefined distribution is unable to model the real distribution [215]. One can simply model the training distribution with histograms [317], [318], [319], [320]. **Kernel density estimation (KDE)** further uses the kernel function as a continuous replacement for the discrete histogram [216], [321], [322]. It flexibly takes parameters such as point weights and bandwidth to control the estimated distribution.

Discussion Although the classic density estimation methods obtain strong AD performance on wide ranges of tasks [323], [324], they are better suited for low-dimensional data. For high-dimensional data in computer vision tasks, these methods suffer from computational and scalability issues due to the curse of dimensionality [325]. To alleviate the problem, shallow methods implement feature engineering to reduce the dimensionality [326], [327].

5.1.2 Density Estimation with Deep Generative Models

In the context of deep learning, neural networks can produce features with high representation quality, which significantly enhance the performance of classic density estimation.

a. AE/VAE-based Models An autoencoder (AE) learns efficient representations of unlabeled data by reconstructing

the input from the latent embedding [328]. Variational autoencoder (VAE) [329] encodes input images into latent vectors under the Gaussian distribution. The learned encoding can be considered as the lower-dimensional representation of the input. Classic density estimation methods can then be applied on top of these deep representations [157], [158], [159].

b. GAN-based Models **Generative adversarial networks (GANs)** consist of a generative network and a discriminative network, contesting with each other in a zero-sum game [330]. Typically, the generative network learns to map from a latent space to a data distribution of interest, whereas the discriminative network distinguishes candidates produced by the generator from the true data distribution. However, unlike the previous AE/VAE paradigm, the lack of an encoder makes it difficult for a GAN to directly find the corresponding embedding for a given image. To solve the problem, ADGAN [160] searches for a good representation in the latent space for a given sample. If such a representation is not found, the sample is deemed anomalous. However, this method can be computationally expensive.

c. Flow-based Models A normalizing flow describes the transformation of a probability density through a sequence of invertible mappings. By repeatedly applying the rule for change of variables, the initial density “flows” through the sequence of invertible mappings [162], [217]. Therefore, methods with the normalizing flow can directly estimate the likelihood of the input space. The flow-based methods are appraised by their elegant mathematical presentations, and are also shown to be sensitive to low-level features only. Flow-based methods can lead to significant computational costs since no dimensionality reduction is performed.

d. Representation Enhancement Apart from obtaining visual embeddings through generative models, some methods focus on enhancing the model capacity to increase the representation power of the extracted features, which may better characterize the normality/ID-ness for more accurate density estimation. These strategies include data augmentation [218], adversarial training [159], distillations [219], loss function enhancement [220], and usage of shallow [221], [222] and local features [223].

5.1.3 Energy-based Models

Energy-based model (EBM) is a generative model that uses a scalar energy score to express the probability density of variables through unnormalized negative log probability [331], which provides a valid solution for AD [224]. However, compared to standard deep learning models, the training process of EBMs can be computationally expensive, since MCMC sampling and approximations are required to calculate integrals. To address the problem, methods such as the score matching method [225] and stochastic gradient Langevin dynamics [226] are proposed for efficient training.

5.1.4 Frequency-based Methods

Previous works also explored frequency domain analysis for anomaly detection. While humans perceive images based on low-frequency components, CNN models can largely depend on high-frequency components for decision-making [227], [228]. Methods such as CNN kernel smoothing [227] and spectrum-oriented data augmentation [229] are proposed to

suppress the influence of high-frequency components. Other works also find that adversarial attacks on low-frequency components are also difficult to detect, therefore proposing methods to target the phase spectrum [230]. Frequency-based methods mainly focus on sensory AD (especially on detecting adversarial examples).

5.2 Reconstruction-based Methods

Similar to Section 4.3, some AD methods rely on the different performance on performing normal and abnormal data to identify anomalies. The difference of model performance can be measured in the feature space (Section 5.2.1) or by the reconstruction error (Section 5.2.2).

5.2.1 Sparse Representation

Sparse reconstruction assumes that every normal sample can be reconstructed accurately using a limited set of basis functions, whereas anomalous data should suffer from larger reconstruction costs, thus generating a dense representation [231], [232], [233]. Exemplar techniques for sparse encoding include L_1 norm-based kernel PCA [234] and low-rank embedded networks [235].

5.2.2 Reconstruction-Error Methods

Reconstruction-error methods rely on the assumption that a reconstruction model trained on the normal data will produce higher-quality outcomes for normal test samples as opposed to anomalies. Deep reconstruction models include AEs [236], VAEs [237], GANs [238], and U-Net [239] that can all be used as the backbone for this method.

a. AE/VAE-based Models Apart from the standard combination of reconstruction-error and AE/VAE models [236], [237], other methods use more sophisticated strategies such as reconstructing by memorized normality [240], [241], adapting model architectures [242], and partial/conditional reconstruction [159], [243], [244]. In the semi-supervised AD setting, CoRA [245] trains two AEs on inliers and outliers respectively. The reconstruction errors derived from the two AEs can be used as an indicator of anomaly.

b. GAN-based Models Advancement in generative modeling has led to the remarkable development of reconstruction-error methods using GANs. The discriminator in GANs intrinsically calculates the reconstruction error for anomaly detection [238]. Moreover, variants of GANs—such as denoising GANs [161] and class-conditional GANs [56]—enable further performance improvement by increasing the reconstruction difficulty. Some methods utilize the performance of the reconstructed image in downstream tasks to further amplify the reconstruction error of anomalies [246]. Ensembling can also enhance the performance [247].

5.2.3 Gradient-based Methods

Gradient-based method belongs to meta-learning or learning to learn, which is a topic of systematically observing the internal mechanisms of the learning tasks or models to propose methods based on the learned experience, or meta-data [332], [333]. To address AD tasks, some method observes the different patterns on training gradient between normalities and anomalies in a reconstruction task and hence use gradient-based representation to characterize anomalies [248].

5.3 Distance-based Methods

Distance-based methods detect anomalies by calculating the distance between targeted samples and a number of internally stored exemplars, or prototypes [334]. These methods usually require training data in the memory. Representative methods include K-nearest Neighbors [249], prototype-based methods [250], [251], as well as methods to be introduced in Section 6.2.

5.4 Classification-based Methods

AD and one-class ND is often formulated as an unsupervised learning problem, where the entire ID data belongs to one class. Fully supervised AD is studied in [335]. The idea of classifier boundaries is successfully implemented and marked as a one-class classification task [252], [336], which we describe in Section 5.4.1. When it comes to semi-supervised AD setting where unlabeled data is introduced for training, PU learning is proposed for this specific problem, which will be introduced in Section 5.4.2. Lastly, we introduce self-supervised learning methods in Section 5.4.3.

5.4.1 One-Class Classification

One-class classification (OCC) directly learns a decision boundary that corresponds to a desired density level set of the normal data distribution [252]. DeepSVDD [253] first introduced the classic OCC to the deep learning community, which maps normal/ID examples into a hypersphere so that the description of normality is bounded. Deviations from this description are then deemed to be anomalies. Later, some works try to extend the method through elastic regularization [254] or constructing an adapted description with multi-linear hyperplanes [255].

5.4.2 Positive-Unlabeled Learning

Positive-unlabeled learning, or PU learning, focuses on the semi-supervised AD setting where unlabeled data is available in addition to the normal data [256], [257], [258]. The unlabeled data can contain both positive and negative examples. Popular PU learning methods generally rely on two strategies. One approach is to select reliable negative samples from unlabeled data and convert them into the supervised AD setting. Techniques such as distance to prototypes [259], [260], [261], clustering [262], [263], and density-based models [264] are used to filter out reliable negatives. Others consider the entire unlabeled set as noisy negatives, converting it into learning with noisy labels setting. Techniques such as sample re-weighting [265] and label cleaning methods [266], [267] have also shown their effectiveness for the task. Besides, reconstruction-error methods can be re-purposed for PU learning by training two reconstruction models for the positive and unlabeled set, respectively [268]. The comparison between their reconstruction-error scores indicates the final AD decision.

5.4.3 Self-Supervised Learning

Self-supervised learning methods tackle the AD and one-class ND problems in two aspects: (1) the enhancement of feature quality can improve AD performance; (2) some well-designed surrogate tasks can help reveal the anomalies from normal samples. In this part, we only discuss the second

pretext task designing, since the first methods that improve feature quality are introduced with their corresponding main tasks, such as in Section 5.1.2. One classic method is isolation forest [269], which generates a random forest to contrast every normal sample. A test anomaly can be isolated in fewer steps than normal instances. Other methods use pretext tasks such as contrastive learning [270] and image transformation prediction [271], [272], where anomalies are more likely to make mistakes on the designed task. For video data, a natural self-supervised task is to predict future frames based on the existing ones [273], where larger error indicates abnormalities.

5.5 Discussion

Sensory vs. Semantic AD Sensory and semantic AD both consider the normality as homogeneous, even though there might be multiple categories in the normal data. Solutions to semantic AD are mostly applicable to sensory AD problems. In particular, sensory AD problems can benefit from methods that focus on lower-level features (*e.g.*, flow-based and hidden feature-based), local representations, and frequency-based methods (*c.f.* Section 5.1.4).

Theoretical Analysis In addition to algorithmic development, several works provided theoretical analysis on AD and one-class ND. In [60], a clean set of ID and a mixed set of ID/OOD are constructed with identical sample sizes. A PAC-style finite sample guarantee is achieved for a certain probability of detecting a certain portion of anomalies with the minimum number of false alarms. Furthermore, in [61], a generalization error bound is provided for PU learning methods in semi-supervised AD.

Anomaly Detection vs. Outlier Detection If we model the test samples and training samples altogether, the AD problem will be transformed into an OD problem, and therefore the transductive approaches in Section 6 are also applicable. However, this method requires all training data to estimate test abnormality, which greatly increases the storage burden and computational complexity. Therefore, we do not include these methods in this part, but leave it to Section 6.

6 OUTLIER DETECTION: METHODOLOGY

Outlier detection (OD) requires the observation of all samples and aims to detect those that deviate significantly from the majority distribution. OD approaches are usually transductive, rather than inductive. Several surveys reviewed methodologies on this topic, yet mostly within the field of data mining [16], [17], [18], [19]. In this section, we briefly review OD methods, especially those developed for computer vision tasks using deep neural networks. We find that although deep learning methods rarely directly solve the OD problem, the data cleaning procedure, which is the prerequisite procedure of learning from open-set noisy data [84], [85] and open set semi-supervised learning [86], are solving the OD tasks.

6.1 Density-based Methods

A basic idea of OD models the entire dataset as a Gaussian distribution and flags samples that are at least three standard

deviations from the mean [337], [338]. Other parametric probabilistic methods make use of Mahalanobis distance [315] and Gaussian mixtures [339] to model the data density. Similar to the “three standard deviations” rules, the interquartile range can also be used to identify outliers [274], forming a classic non-parametric probabilistic method. Local outlier factor (LOF) estimate the density of a given point via the ratio of the local reachability of its neighbors and itself [275], followed by derivations for robustification [276], [277] and simplification [278]. RANSAC [279] iteratively estimates parameters of a mathematical model to fit the data and find the samples as outliers that contribute less to estimates. Generally, classic density methods for AD such as kernel density estimation (*c.f.* Section 5.1) are also applicable for OD. Although these methods suffer from the curse of dimensionality on images, they can be alleviated by dimensionality reduction methods [280], [281] and the NN-based density methods (*c.f.* Section 5.1).

6.2 Distance-based Methods

A simple method to detect outliers is counting the number of neighbors within a certain radius, or measure the k -th-nearest neighbor distance [340], [341]. We mainly discuss cluster-based methods and graph-based methods here.

6.2.1 Cluster-based Methods

DBSCAN [282] accumulates samples based on the distance-based density to form clusters. Samples that lie outside the major clusters are recognized as outliers. Subsequent works improve the clustering approaches by considering the confidence of cluster labels [283].

6.2.2 Graph-based Methods

Another set of methods uses the relationship among data points and constructs a neighborhood graph [342], [343] or its variants [344]. Graph properties and graph mining techniques are employed to find abnormal samples [284], [285], such as graph-based clustering [286], [287], partitioning [345], and label propagation with graph neural networks [288].

6.3 Classification-based Methods

AD methods (*e.g.*, Isolation Forest [269], OC-SVM [252], [253], *etc.*) are also applicable to OD setting. When there are multiple classes in the dataset, researchers find that deep learning models—when trained with outliers—can still show robust prediction capability and identify the outliers [289]. Data cleaning using large pre-trained models is also common in the industry. Techniques to enhance model robustness and feature generalizability can be useful for this task, such as ensembling [290], co-training [291], and distillation [289], [292].

6.4 Discussion

Although the application of OD is not as common as other sub-tasks in the computer vision community, techniques for OD can be valuable for other tasks such as open-set semi-supervised learning, learning with open-set noisy labels, and potentially novelty discovery. In addition to methods reviewed here, most solutions to AD/ND/OOD detection

can also be applied by considering all observations as ID (for model training) and then applying the model again on all the observations. In this case, methods such as reconstruction-based PCA (*c.f.* Section 5.2) and energy-based models (*c.f.* Section 5.1.3 and 3.1.1) can also address the OD problem.

7 CHALLENGES AND FUTURE DIRECTIONS

In this section, we discuss the challenges and future directions of **generalized OOD detection**.

7.1 Challenges

a. Proper Evaluation and Benchmarking We hope this survey can clarify the distinctions and connections of various sub-tasks, and help future works properly identify the target problem and benchmarks within the framework. The mainstream OOD detection works primarily focus on detecting semantic shifts. Admittedly, the field of OOD detection can be very broad due to the **diverse nature of distribution shifts**. Such a broad OOD definition also leads to some challenges and concerns [35], [141], which advocate a clear specification of OOD type in consideration (*e.g.*, **semantic OOD**, **adversarial OOD**, *etc.*) so that proposed solutions can be more specialized. Besides, the motivation of **detecting a certain distribution shift also requires clarification**. While rejecting classifying samples with semantic shift is apparent, detecting sensory OOD should be specified to some meaningful scenarios to contextualize the necessity and practical relevance of the task.

We also urge the community to carefully **construct the benchmarks and evaluations**. It is noticed that early work [64] ignored the fact that **some OOD datasets may contain images with ID categories**, causing inaccurate performance evaluation. Fortunately, recent OOD detection works [131] have realized this flaw and **pay special attention to removing ID classes from OOD samples to ensure proper evaluation**.

b. Outlier-free OOD Detection The **outlier exposure** approach [67] imposes a **strong assumption** of the availability of OOD training data, which can be difficult to obtain in practice. Moreover, one needs to perform careful deduplication to ensure that the outlier training data does not contain ID data. These restrictions may lead to inflexible solutions and prevent the adoption of methods in the real world. Going forward, a major challenge for the field is to devise outlier-free learning objectives that are less dependent on auxiliary outlier dataset.

c. Tradeoff Between Classification and OOD Detection In OSR and OOD detection, it is important to achieve the dual objectives simultaneously: one for the ID task (*e.g.*, image classification), another for the OOD detection task. For a shared network, an inherent trade-off may exist between the two tasks. Promising solutions should strive for both. These two tasks may or may not contradict each other, depending on the methodologies. For example, [100] advocated the integration of image classification and open-set recognition so that the model will possess the capability of discriminative recognition on known classes and sensitivity to novel classes at the same time. [314] also showed that the ability of detecting novel classes can be highly correlated with its accuracy on the closed-set classes. [131] demonstrated that optimizing

for the cluster compactness of ID classes may facilitate both improved classification and distance-based OOD detection performance. Such solutions may be more desirable than ND, which develops a binary OOD detector separately from the classification model, and requires deploying two models.

d. Real-world Benchmarks and Evaluations Current methods have been primarily evaluated on **small data sets** such as CIFAR. It's been shown that approaches developed on the **CIFAR benchmark** might **not translate effectively** into ImageNet benchmark with a large semantic space, highlighting the need to evaluate OOD detection in a large-scale real-world setting. Therefore, we encourage future research to evaluate on ImageNet-based OOD detection benchmark [138], as well as large-scale OSR benchmark [314], and test the limits of the method developed. Moreover, real-world benchmarks that go beyond image classification can be valuable for the research community. In particular, recent works start looking more into object-level OOD detection [146], [147], which can be useful for safety-critical settings such as autonomous driving.

7.2 Future Directions

a. Methodologies across Sub-tasks Due to the inherent connections among different sub-tasks, their solution space can be shared and inspired from each other. For example, the recent emerging **density-based** OOD detection research (*c.f.* Section 3.2) can draw insights from the density-based AD methods (*c.f.* Section 5.1) that have been around for a long time.

b. OOD Detection & Generalization An open-world classifier should consider two tasks, *i.e.*, **being robust to covariate shift while being aware of the semantic shift**. Existing works pursue these two goals independently. Recent work proposes a semantically coherent OOD detection framework [131] that encourages detecting semantic OOD samples while being robust to negligible covariate shift. Given the vague definition of OOD, [346] proposed a formalization of OOD detection by explicitly taking into account the separation between invariant features (semantic related) and environmental features (non-semantic). The work highlighted that spurious environmental features in the training set can significantly impact OOD detection, especially when the semantic OOD data contains the spurious feature. Further, full-spectrum OOD detection [308] highlights the effects of “covariate-shifted in-distribution”, and show that most of the previous OOD detectors are unfortunately sensitive to covariate shift rather than semantic shift. This setting explicitly promotes the generalization ability of OOD detectors. Recent works on open long-tailed recognition [100], open compound domain adaptation [92], open-set domain adaptation [347] and open-set domain generalization [348] consider the potential existence of open-class samples. Looking ahead, we envision great research opportunities on how OOD detection and OOD generalization can better enable each other [100], in terms of both algorithmic design and comprehensive performance evaluation.

c. OOD Detection & Open-Set Noisy Labels Existing methods of learning from open-set noisy labels focus on **suppressing the negative effects of noise** [84], [349]. However, the open-set noisy samples can be useful for outlier exposure

(c.f. Section 3.1.1) [345] and potentially benefit OOD detection. With a similar idea, the setting of open-set semi-supervised learning can be promising for OOD detection. We believe the combination between OOD detection and the previous two fields can provide more insights and possibilities.

d. Theoretical Analysis While most of the existing OOD detection works focus on developing effective approaches to obtain better empirical performance, the theoretical analysis remains largely underexplored. Recently, [293] developed an analytical framework that provided theoretical understanding for OOD detection. We hope future research can also contribute theoretical analyses and provide in-depth insights that help guide algorithmic development with rigorous guarantees.

e. OOD Detection For Broader Learning Tasks As mentioned in the Section 3.5, OOD detection encompasses a broader spectrum of learning tasks, including multi-label classification [106], object detection [146], [147], image segmentation [68], time-series prediction [350], and LiDAR-based 3D object detection [351]. For the classification task itself, the researchers also extended the OOD detection technique to improve the reliability of zero-shot pretrained models [352] (e.g., CLIP). Other works explore how to use OOD detection methods to only generate reliable image captions [353]. We believe that with OOD detection methods, many more fields can benefit from reliable and practical models. In turn, insights from various fields can also help OOD detection to develop better.

8 CONCLUSION

In this survey, we comprehensively review five topics: AD, ND, OSR, OOD detection, and OD, and unify them as a framework of *generalized OOD detection*. By articulating the motivations and definitions of each sub-task, we encourage follow-up works to accurately locate their target problems and find the most suitable benchmarks. By sorting out the methodologies for each sub-task, we hope that readers can easily grasp the mainstream methods, identify suitable baselines, and contribute future solutions in light of existing ones. By providing insights, challenges, and future directions, we hope that future works will pay more attention to the existing problems and explore more interactions across other tasks within or even outside the scope of generalized OOD detection.

ACKNOWLEDGMENTS

This study is supported by NTU NAP, and the RIE2020 Industry Alignment Fund – Industry Collaboration Projects (IAF-ICP) Funding Initiative, as well as cash and in-kind contributions from the industry partner(s). YL is supported by the Office of the Vice Chancellor for Research and Graduate Education (OVCRG) with funding from the Wisconsin Alumni Research Foundation (WARF).

REFERENCES

- [1] D. Amodei, C. Olah, J. Steinhardt, P. Christiano, J. Schulman, and D. Mané, “Concrete problems in AI safety,” *arXiv preprint arXiv:1606.06565*, 2016. 1
- [2] T. G. Dietterich, “Steps toward robust artificial intelligence,” *AI Magazine*, 2017. 1
- [3] J. Leike, M. Martic, V. Krakovna, P. A. Ortega, T. Everitt, A. Lefrancq, L. Orseau, and S. Legg, “Ai safety gridworlds,” *arXiv preprint arXiv:1711.09883*, 2017. 1
- [4] N. A. Smuha, “The EU approach to ethics guidelines for trustworthy artificial intelligence,” *Computer Law Review International*, 2019. 1
- [5] B. Shneiderman, “Bridging the gap between ethics and practice: Guidelines for reliable, safe, and trustworthy human-centered ai systems,” *TiiS*, 2020. 1
- [6] S. Mohseni, H. Wang, Z. Yu, C. Xiao, Z. Wang, and J. Yadawa, “Practical machine learning safety: A survey and primer,” *arXiv preprint arXiv:2106.04823*, 2021. 1
- [7] D. Hendrycks, N. Carlini, J. Schulman, and J. Steinhardt, “Unsolved problems in ml safety,” *arXiv preprint arXiv:2109.13916*, 2021. 1
- [8] D. Hendrycks and M. Mazeika, “X-risk analysis for ai research,” *arXiv preprint arXiv:2206.05862*, 2022. 1
- [9] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *NIPS*, 2012. 1
- [10] K. He, X. Zhang, S. Ren, and J. Sun, “Delving deep into rectifiers: Surpassing human-level performance on imagenet classification,” in *ICCV*, 2015. 1
- [11] N. Drummond and R. Shearer, “The open world assumption,” in *eSI Workshop*, 2006. 1
- [12] D. Hendrycks and K. Gimpel, “A baseline for detecting misclassified and out-of-distribution examples in neural networks,” in *ICLR*, 2017. 1
- [13] S. Ben-David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, and J. W. Vaughan, “A theory of learning from different domains,” *Machine learning*, 2010. 1
- [14] D. Li, Y. Yang, Y.-Z. Song, and T. M. Hospedales, “Deeper, broader and artier domain generalization,” in *ICCV*, 2017. 1
- [15] M. Wang and W. Deng, “Deep visual domain adaptation: A survey,” *Neurocomputing*, 2018. 1, 6
- [16] C. C. Aggarwal and P. S. Yu, “Outlier detection for high dimensional data,” in *ACM SIGMOD*, 2001. 2, 5, 14
- [17] V. Hodge and J. Austin, “A survey of outlier detection methodologies,” *Artificial intelligence review*, 2004. 2, 5, 14
- [18] I. Ben-Gal, “Outlier detection,” in *Data mining and knowledge discovery handbook*, 2005. 2, 5, 14
- [19] H. Wang, M. J. Bah, and M. Hammad, “Progress in outlier detection techniques: A survey,” *Ieee Access*, 2019. 2, 14
- [20] L. Ruff, J. R. Kauffmann, R. A. Vandermeulen, G. Montavon, W. Samek, M. Kloft, T. G. Dietterich, and K.-R. Müller, “A unifying review of deep and shallow anomaly detection,” *Proceedings of the IEEE*, 2021. 2, 3, 4, 6
- [21] G. Pang, C. Shen, L. Cao, and A. v. d. Hengel, “Deep learning for anomaly detection: A review,” *arXiv preprint arXiv:2007.02500*, 2020. 2, 4, 6
- [22] S. Bulusu, B. Kailkhura, B. Li, P. K. Varshney, and D. Song, “Anomalous example detection in deep learning: A survey,” *IEEE Access*, 2020. 2, 6
- [23] R. Chalapathy and S. Chawla, “Deep learning for anomaly detection: A survey,” *arXiv preprint arXiv:1901.03407*, 2019. 2, 4, 6
- [24] M. A. Pimentel, D. A. Clifton, L. Clifton, and L. Tarassenko, “A review of novelty detection,” *Signal Processing*, 2014. 2
- [25] D. Miljković, “Review of novelty detection methods,” in *MIPRO*, 2010. 2
- [26] M. Markou and S. Singh, “Novelty detection: a review—part 1: statistical approaches,” *Signal processing*, 2003. 2, 4
- [27] M. Markou and S. Singh, “Novelty detection: a review—part 2: neural network based approaches,” *Signal processing*, 2003. 2, 4
- [28] T. E. Boult, S. Cruz, A. R. Dhamija, M. Gunther, J. Henrydoss, and W. J. Scheirer, “Learning and the unknown: Surveying steps toward open world recognition,” in *AAAI*, 2019. 2, 5
- [29] C. Geng, S.-j. Huang, and S. Chen, “Recent advances in open set recognition: A survey,” *TPAMI*, 2020. 2
- [30] A. Mahdavi and M. Carvalho, “A survey on open set recognition,” *arXiv preprint arXiv:2109.00893*, 2021. 2
- [31] I. J. Goodfellow, J. Shlens, and C. Szegedy, “Explaining and harnessing adversarial examples,” in *ICLR*, 2015. 3
- [32] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, “Towards deep learning models resistant to adversarial attacks,” *ICLR*, 2018. 3

- [33] J. Quiñonero-Candela, M. Sugiyama, N. D. Lawrence, and A. Schwaighofer, *Dataset shift in machine learning*. Mit Press, 2009. 3
- [34] L. A. Gatys, A. S. Ecker, and M. Bethge, "Image style transfer using convolutional neural networks," in *CVPR*, 2016. 3
- [35] F. Ahmed and A. Courville, "Detecting semantic anomalies," in *AAAI*, 2020. 3, 15
- [36] L. Zhang, M. Goldstein, and R. Ranganath, "Understanding failures in out-of-distribution detection with deep generative models," in *ICML*, 2021. 3
- [37] N. Akhtar and A. Mian, "Threat of adversarial attacks on deep learning in computer vision: A survey," *IEEE Access*, 2018. 3
- [38] K. Patel, H. Han, and A. K. Jain, "Secure face unlock: Spoof detection on smartphones," *IEEE transactions on information forensics and security*, 2016. 3
- [39] D. Wen, H. Han, and A. K. Jain, "Face spoof detection with image distortion analysis," *IEEE Transactions on Information Forensics and Security*, 2015. 3
- [40] K. A. Nixon, V. Aimale, and R. K. Rowe, "Spoof detection schemes," in *Handbook of biometrics*, 2008. 3
- [41] G. Polatkan, S. Jafarpour, A. Brasoveanu, S. Hughes, and I. Daubechies, "Detection of forgery in paintings using supervised learning," in *ICIP*, 2009. 3
- [42] B. Dolhansky, R. Howes, B. Pflaum, N. Baram, and C. C. Ferrer, "The deepfake detection challenge (dfdc) preview dataset," *arXiv preprint arXiv:1910.08854*, 2019. 3
- [43] L. Jiang, Z. Guo, W. Wu, Z. Liu, Z. Liu, C. C. Loy, S. Yang, Y. Xiong, W. Xia, B. Chen, P. Zhuang, S. Li, S. Chen, T. Yao, S. Ding, J. Li, F. Huang, L. Cao, R. Ji, C. Lu, and G. Tan, "DeeperForensics Challenge 2020 on real-world face forgery detection: Methods and results," *arXiv preprint arXiv:2102.09471*, 2021. 3
- [44] P. Yang, D. Baracchi, R. Ni, Y. Zhao, F. Argenti, and A. Piva, "A survey of deep learning-based source image forensics," *Journal of Imaging*, 2020. 3
- [45] P. Bergmann, M. Fauser, D. Sattlegger, and C. Steger, "Mvtac ad-a comprehensive real-world dataset for unsupervised anomaly detection," in *CVPR*, 2019. 3
- [46] W.-H. Chu and K. M. Kitani, "Neural batch sampling with reinforcement learning for semi-supervised anomaly detection," in *ECCV*, 2020. 3
- [47] D. J. Atha and M. R. Jahanshahi, "Evaluation of deep learning approaches based on convolutional neural networks for corrosion detection," *Structural Health Monitoring*, 2018. 3
- [48] H. Idrees, M. Shah, and R. Surette, "Enhancing camera surveillance using computer vision: a research note," *Policing: An International Journal*, 2018. 3, 4
- [49] C. P. Diehl and J. B. Hampshire, "Real-time object classification and novelty detection for collaborative video surveillance," in *IJCNN*, 2002. 3, 4
- [50] L.-J. Li and L. Fei-Fei, "Optimol: automatic online picture collection via incremental model learning," *IJCV*, 2010. 3
- [51] T. Fawcett, "An introduction to roc analysis," *Pattern recognition letters*, 2006. 4
- [52] D. M. Powers, "Evaluation: from precision, recall and f-measure to roc, informedness, markedness and correlation," *JMLT*, 2020. 4
- [53] H. R. Kerner, D. F. Wellington, K. L. Wagstaff, J. F. Bell, C. Kwan, and H. B. Amor, "Novelty detection for multispectral images with application to planetary exploration," in *AAAI*, 2019. 4
- [54] H. Al-Behadili, A. Grumpe, and C. Wöhler, "Incremental learning and novelty detection of gestures in a multi-class system," in *AIMS*, 2015. 4
- [55] D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell, "Curiosity-driven exploration by self-supervised prediction," in *ICML*, 2017. 4
- [56] P. Perera, R. Nallapati, and B. Xiang, "Ocgan: One-class novelty detection using gans with constrained latent representations," in *CVPR*, 2019. 4, 7, 13
- [57] Y. Xia, X. Cao, F. Wen, G. Hua, and J. Sun, "Learning discriminative reconstructions for unsupervised outlier removal," in *CVPR*, 2015. 4
- [58] W. J. Scheirer, A. de Rezende Rocha, A. Sapkota, and T. E. Boult, "Toward open set recognition," *TPAMI*, 2013. 4, 10, 11
- [59] A. R. Dhamija, M. Günther, and T. E. Boult, "Reducing network agnostophobia," in *NeurIPS*, 2018. 4, 5, 7, 8
- [60] S. Liu, R. Garrepalli, T. Dietterich, A. Fern, and D. Hendrycks, "Open category detection with pac guarantees," in *ICML*, 2018. 4, 7, 14
- [61] Z. Fang, J. Lu, A. Liu, F. Liu, and G. Zhang, "Learning bounds for open-set learning," in *ICML*, 2021. 4, 7, 14
- [62] E. Sorio, A. Bartoli, G. Davanzo, and E. Medvet, "Open world classification of printed invoices," in *Proceedings of the 10th ACM symposium on Document engineering*, 2010. 4
- [63] H. Xu, B. Liu, L. Shu, and P. Yu, "Open-world learning and application to product classification," in *WWW*, 2019. 4
- [64] D. Hendrycks and K. Gimpel, "A baseline for detecting misclassified and out-of-distribution examples in neural networks," in *ICLR*, 2017. 5, 7, 15
- [65] X. Huang, D. Kroening, W. Ruan, J. Sharp, Y. Sun, E. Thamo, M. Wu, and X. Yi, "A survey of safety and trustworthiness of deep neural networks: Verification, testing, adversarial attack and defence, and interpretability," *Computer Science Review*, 2020. 5
- [66] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *CVPR*, 2012. 5
- [67] D. Hendrycks, M. Mazeika, and T. Dietterich, "Deep anomaly detection with outlier exposure," in *ICLR*, 2019. 5, 7, 8, 15
- [68] D. Hendrycks, S. Basart, M. Mazeika, M. Mostajabi, J. Steinhardt, and D. Song, "Scaling out-of-distribution detection for real-world settings," in *ICML*, 2022. 5, 16
- [69] Y.-C. Hsu, Y. Shen, H. Jin, and Z. Kira, "Generalized odin: Detecting out-of-distribution image without learning from out-of-distribution data," in *CVPR*, 2020. 5, 7, 8
- [70] G. Pleiss, A. Souza, J. Kim, B. Li, and K. Q. Weinberger, "Neural network out-of-distribution detection for regression tasks," 2019. 5
- [71] O. Vinyals, T. Ewalds, S. Bartunov, P. Georgiev, A. S. Vezhnevets, M. Yeo, A. Makhzani, H. Küttler, J. Agapiou, J. Schriftwieser, et al., "Starcraft ii: A new challenge for reinforcement learning," *arXiv preprint arXiv:1708.04782*, 2017. 5
- [72] A. Sedlmeier, T. Gabor, T. Phan, L. Belzner, and C. Linnhoff-Popien, "Uncertainty-based out-of-distribution detection in deep reinforcement learning," *arXiv preprint arXiv:1901.02219*, 2019. 5
- [73] D. Zimmerer, P. M. Full, F. Isensee, P. Jäger, T. Adler, J. Petersen, G. Köhler, T. Ross, A. Reinke, A. Kascenas, et al., "Mood 2020: A public benchmark for out-of-distribution detection and localization on medical images," *IEEE Transactions on Medical Imaging*, 2022. 5
- [74] M. I. Tariq, N. A. Memon, S. Ahmed, S. Tayyaba, M. T. Mushtaq, N. A. Mian, M. Imran, and M. W. Ashraf, "A review of deep learning security and privacy defensive techniques," *Mobile Information Systems*, vol. 2020, 2020. 5
- [75] Wikipedia contributors, "Outlier from Wikipedia, the free encyclopedia," 2021. [Online; accessed 12 August 2021]. 5
- [76] M. Bianchini, A. Belahcen, and F. Scarselli, "A comparative study of inductive and transductive learning with feedforward neural networks," in *Conference of the Italian Association for Artificial Intelligence*, 2016. 5
- [77] I. Ben-Gal, "Outlier detection," in *Data mining and knowledge discovery handbook*, 2005. 6
- [78] S. Basu and M. Meckesheimer, "Automatic outlier detection for time series: an application to sensor data," *Knowledge and Information Systems*, 2007. 6
- [79] Y. Dou, W. Li, Z. Liu, Z. Dong, J. Luo, and S. Y. Philip, "Uncovering download fraud activities in mobile app markets," in *ASONAM*, 2019. 6
- [80] T. Xiao, C. Zhang, and H. Zha, "Learning to detect anomalies in surveillance video," *IEEE Signal Processing Letters*, 2015. 6
- [81] H. Liu, S. Shah, and W. Jiang, "On-line outlier detection and data cleaning," *Computers & chemical engineering*, 2004. 6
- [82] A. Loureiro, L. Torgo, and C. Soares, "Outlier detection using clustering methods: a data cleaning application," in *Proceedings of KDNet Symposium on Knowledge-based Systems*, 2004. 6
- [83] J. Van den Broeck, S. Argeseanu Cunningham, R. Eeckels, and K. Herbst, "Data cleaning: detecting, diagnosing, and editing data abnormalities," *PLoS medicine*, 2005. 6
- [84] Y. Wang, W. Liu, X. Ma, J. Bailey, H. Zha, L. Song, and S.-T. Xia, "Iterative learning with open-set noisy labels," in *CVPR*, 2018. 6, 14, 15
- [85] X. Chen and A. Gupta, "Weby learning of convolutional networks," in *ICCV*, 2015. 6, 14
- [86] K. Cao, M. Brbic, and J. Leskovec, "Open-world semi-supervised learning," *arXiv preprint arXiv:2102.03526*, 2021. 6, 14
- [87] C. Chow, "On optimum recognition error and reject tradeoff," *IEEE Transactions on Information Theory*, 1970. 6

- [88] G. Fumera and F. Roli, "Support vector machines with embedded reject option," in *International Workshop on Support Vector Machines*, 2002. [6](#)
- [89] C. Cortes and V. Vapnik, "Support-vector networks," *Machine learning*, 1995. [6](#)
- [90] A. Nguyen, J. Yosinski, and J. Clune, "Deep neural networks are easily fooled: High confidence predictions for unrecognizable images," in *CVPR*, 2015. [6](#)
- [91] K. Zhou, Z. Liu, Y. Qiao, T. Xiang, and C. C. Loy, "Domain generalization: A survey," *arXiv preprint arXiv:2103.02503*, 2021. [6](#)
- [92] Z. Liu, Z. Miao, X. Pan, X. Zhan, D. Lin, S. X. Yu, and B. Gong, "Open compound domain adaptation," in *CVPR*, 2020. [6, 15](#)
- [93] K. Han, A. Vedaldi, and A. Zisserman, "Learning to discover novel visual categories via deep transfer clustering," in *CVPR*, 2019. [6](#)
- [94] B. Zhao and K. Han, "Novel visual category discovery with dual ranking statistics and mutual knowledge distillation," *NeurIPS*, 2021. [6](#)
- [95] X. Jia, K. Han, Y. Zhu, and B. Green, "Joint representation learning and novel category discovery on single-and multi-modal data," in *ICCV*, 2021. [6](#)
- [96] S. Vaze, K. Han, A. Vedaldi, and A. Zisserman, "Generalized category discovery," in *CVPR*, 2022. [6](#)
- [97] K. Joseph, S. Paul, G. Aggarwal, S. Biswas, P. Rai, K. Han, and V. N. Balasubramanian, "Novel class discovery without forgetting," in *ECCV*, 2022. [6](#)
- [98] W. Wang, V. W. Zheng, H. Yu, and C. Miao, "A survey of zero-shot learning: Settings, methods, and applications," *TIST*, 2019. [6](#)
- [99] A. Bendale and T. Boult, "Towards open world recognition," in *CVPR*, 2015. [6](#)
- [100] Z. Liu, Z. Miao, X. Zhan, J. Wang, B. Gong, and S. X. Yu, "Large-scale long-tailed recognition in an open world," in *CVPR*, 2019. [6, 15](#)
- [101] M. Salehi, H. Mirzaei, D. Hendrycks, Y. Li, M. H. Rohban, and M. Sabokrou, "A unified survey on anomaly, novelty, open-set, and out-of-distribution detection: Solutions and future challenges," *arXiv preprint arXiv:2110.14051*, 2021. [6](#)
- [102] S. Liang, Y. Li, and R. Srikant, "Enhancing the reliability of out-of-distribution image detection in neural networks," in *ICLR*, 2018. [7, 8, 9](#)
- [103] K. Lee, K. Lee, H. Lee, and J. Shin, "A simple unified framework for detecting out-of-distribution samples and adversarial attacks," in *NeurIPS*, 2018. [7, 8, 10](#)
- [104] W. Liu, X. Wang, J. D. Owens, and Y. Li, "Energy-based out-of-distribution detection," in *NeurIPS*, 2020. [7, 8](#)
- [105] C. S. Sastry and S. Oore, "Detecting out-of-distribution examples with gram matrices," in *ICML*, 2020. [7, 8](#)
- [106] H. Wang, W. Liu, A. Bocchieri, and Y. Li, "Can multi-label classification networks know what they don't know?," *NeurIPS*, 2021. [7, 8, 10, 16](#)
- [107] Y. Sun, C. Guo, and Y. Li, "React: Out-of-distribution detection with rectified activations," in *NeurIPS*, 2021. [7, 8](#)
- [108] X. Dong, J. Guo, A. Li, W.-T. Ting, C. Liu, and H. Kung, "Neural mean discrepancy for efficient out-of-distribution detection," in *CVPR*, 2022. [7, 8](#)
- [109] Y. Sun and Y. Li, "Dice: Leveraging sparsification for out-of-distribution detection," in *ECCV*, 2022. [7, 8](#)
- [110] Y. Sun, Y. Ming, X. Zhu, and Y. Li, "Out-of-distribution detection with deep nearest neighbors," in *ICML*, 2022. [7, 10](#)
- [111] Z. Lin, S. D. Roy, and Y. Li, "Mood: Multi-level out-of-distribution detection," in *CVPR*, 2021. [7, 8](#)
- [112] T. DeVries and G. W. Taylor, "Learning confidence for out-of-distribution detection in neural networks," *arXiv preprint arXiv:1802.04865*, 2018. [7, 8](#)
- [113] Y. Wang, B. Li, T. Che, K. Zhou, Z. Liu, and D. Li, "Energy-based open-world uncertainty modeling for confidence calibration," in *ICCV*, 2021. [7, 8](#)
- [114] M. Hein, M. Andriushchenko, and J. Bitterwolf, "Why relu networks yield high-confidence predictions far away from the training data and how to mitigate the problem," in *CVPR*, 2019. [7, 8](#)
- [115] A. Vyas, N. Jammalamadaka, X. Zhu, D. Das, B. Kaul, and T. L. Willke, "Out-of-distribution detection using an ensemble of self supervised leave-out classifiers," in *ECCV*, 2018. [7, 8](#)
- [116] J. Bitterwolf, A. Meinke, and M. Hein, "Certifiably adversarially robust detection of out-of-distribution data," in *NeurIPS*, 2020. [7, 8](#)
- [117] J. Chen, Y. Li, X. Wu, Y. Liang, and S. Jha, "Robust out-of-distribution detection for neural networks," *arXiv preprint arXiv:2003.09711*, 2020. [7, 8](#)
- [118] M. Hein, M. Andriushchenko, and J. Bitterwolf, "Why relu networks yield high-confidence predictions far away from the training data and how to mitigate the problem," in *CVPR*, 2019. [7, 8](#)
- [119] S. Choi and S.-Y. Chung, "Novelty detection via blurring," in *ICLR*, 2020. [7, 8](#)
- [120] S. Thulasidasan, G. Chennupati, J. Bilmes, T. Bhattacharya, and S. Michalak, "On mixup training: Improved calibration and predictive uncertainty for deep neural networks," in *NeurIPS*, 2019. [7, 8](#)
- [121] S. Yun, D. Han, S. J. Oh, S. Chun, J. Choe, and Y. Yoo, "Cutmix: Regularization strategy to train strong classifiers with localizable features," in *CVPR*, 2019. [7, 8](#)
- [122] T. DeVries and G. W. Taylor, "Improved regularization of convolutional neural networks with cutout," *arXiv preprint arXiv:1708.04552*, 2017. [7, 8](#)
- [123] D. Hendrycks, N. Mu, E. D. Cubuk, B. Zoph, J. Gilmer, and B. Lakshminarayanan, "Augmix: A simple data processing method to improve robustness and uncertainty," *arXiv preprint arXiv:1912.02781*, 2019. [7, 8](#)
- [124] A. Meinke and M. Hein, "Towards neural networks that provably know when they don't know," *arXiv preprint arXiv:1909.12180*, 2019. [7, 8](#)
- [125] H. Wei, R. Xie, H. Cheng, L. Feng, B. An, and Y. Li, "Mitigating neural network overconfidence with logit normalization," in *ICML*, 2022. [7, 8](#)
- [126] Q. Yu and K. Aizawa, "Unsupervised out-of-distribution detection by maximum classifier discrepancy," in *ICCV*, 2019. [7, 8](#)
- [127] A.-A. Papadopoulos, M. R. Rajati, N. Shaikh, and J. Wang, "Outlier exposure with confidence control for out-of-distribution detection," *Neurocomputing*, 2021. [7, 8](#)
- [128] Y. Li and N. Vasconcelos, "Background data resampling for outlier-aware classification," in *CVPR*, 2020. [7, 8](#)
- [129] J. Chen, Y. Li, X. Wu, Y. Liang, and S. Jha, "Atom: Robustifying out-of-distribution detection using outlier mining," *ECML&PKDD*, 2021. [7, 8](#)
- [130] S. Mohseni, M. Pitale, J. Yadawa, and Z. Wang, "Self-supervised learning for generalizable out-of-distribution detection," in *AAAI*, 2020. [7, 8](#)
- [131] J. Yang, H. Wang, L. Feng, X. Yan, H. Zheng, W. Zhang, and Z. Liu, "Semantically coherent out-of-distribution detection," in *ICCV*, 2021. [7, 9, 15](#)
- [132] S. Thulasidasan, S. Thapa, S. Dhaubhadel, G. Chennupati, T. Bhattacharya, and J. Bilmes, "An effective baseline for robustness to distributional shift," *arXiv preprint arXiv:2105.07107*, 2021. [7, 8](#)
- [133] A. Shafaei, M. Schmidt, and J. J. Little, "A less biased evaluation of out-of-distribution sample detectors," in *BMVC*, 2019. [7, 9](#)
- [134] S. Mohseni, M. Pitale, J. Yadawa, and Z. Wang, "Self-supervised learning for generalizable out-of-distribution detection," in *AAAI*, 2020. [7, 8](#)
- [135] J. Katz-Samuels, J. Nakhleh, R. Nowak, and Y. Li, "Training ood detectors in their natural habitats," in *International Conference on Machine Learning (ICML)*, PMLR, 2022. [7, 9](#)
- [136] Y. Ming, Y. Fan, and Y. Li, "Poem: Out-of-distribution detection with posterior sampling," in *ICML*, 2022. [7, 8](#)
- [137] K. Lee, K. Lee, K. Min, Y. Zhang, J. Shin, and H. Lee, "Hierarchical novelty detection for visual object recognition," in *CVPR*, 2018. [7, 9](#)
- [138] R. Huang and Y. Li, "Mos: Towards scaling out-of-distribution detection for large semantic space," in *CVPR*, 2021. [7, 9, 15](#)
- [139] G. Shalev, Y. Adi, and J. Keshet, "Out-of-distribution detection using multiple semantic label representations," in *NeurIPS*, 2018. [7, 9](#)
- [140] S. Fort, J. Ren, and B. Lakshminarayanan, "Exploring the limits of out-of-distribution detection," *NeurIPS*, 2021. [7, 9](#)
- [141] W. Gan, "Language guided out-of-distribution detection," 2021. [7, 9, 15](#)
- [142] K. Lee, H. Lee, K. Lee, and J. Shin, "Training confidence-calibrated classifiers for detecting out-of-distribution samples," 2018. [7, 9](#)
- [143] S. Vernekar, A. Gaurav, V. Abdelzad, T. Denouden, R. Salay, and K. Czarnecki, "Out-of-distribution detection in classifiers via generation," in *NeurIPS-W*, 2019. [7, 9](#)

- [144] K. Sricharan and A. Srivastava, "Building robust classifiers through generation of confident out-of-distribution examples," in *NeurIPS-W*, 2018. [7](#) [9](#)
- [145] D. Mandal, S. Narayan, S. K. Dwivedi, V. Gupta, S. Ahmed, F. S. Khan, and L. Shao, "Out-of-distribution detection for generalized zero-shot action recognition," in *CVPR*, 2019. [7](#)
- [146] X. Du, Z. Wang, M. Cai, and Y. Li, "Vos: Learning what you don't know by virtual outlier synthesis," in *Proceedings of the International Conference on Learning Representations*, 2022. [7](#) [9](#), [10](#), [15](#), [16](#)
- [147] X. Du, X. Wang, G. Gozum, and Y. Li, "Unknown-aware object detection: Learning what you don't know from videos in the wild," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. [7](#) [9](#), [10](#), [15](#), [16](#)
- [148] R. Huang, A. Geng, and Y. Li, "On the importance of gradients for detecting distributional shifts in the wild," in *NeurIPS*, 2021. [7](#) [9](#)
- [149] Y. Gal and Z. Ghahramani, "Dropout as a bayesian approximation: Representing model uncertainty in deep learning," in *ICML*, 2016. [7](#) [9](#)
- [150] B. Lakshminarayanan, A. Pritzel, and C. Blundell, "Simple and scalable predictive uncertainty estimation using deep ensembles," *NeurIPS*, 2017. [7](#) [9](#)
- [151] K. Osawa, S. Swaroop, A. Jain, R. Eschenhagen, R. E. Turner, R. Yokota, and M. E. Khan, "Practical deep learning with bayesian principles," in *NeurIPS*, 2019. [7](#) [9](#)
- [152] A. Malinin and M. Gales, "Predictive uncertainty estimation via prior networks," in *NeurIPS*, 2018. [7](#) [9](#)
- [153] A. Malinin and M. Gales, "Reverse kl-divergence training of prior networks: Improved uncertainty and adversarial robustness," in *NeurIPS*, 2019. [7](#) [9](#)
- [154] J. Nandy, W. Hsu, and M. L. Lee, "Towards maximizing the representation gap between in-domain & out-of-distribution examples," in *NeurIPS*, 2020. [7](#) [9](#)
- [155] D. Hendrycks, X. Liu, E. Wallace, A. Dziedzic, R. Krishnan, and D. Song, "Pretrained transformers improve out-of-distribution robustness," *arXiv preprint arXiv:2004.06100*, 2020. [7](#) [9](#)
- [156] R. Koner, P. Sinhamahapatra, K. Roscher, S. Günnemann, and V. Tresp, "Oodformer: Out-of-distribution detection transformer," *arXiv preprint arXiv:2107.08976*, 2021. [7](#) [9](#)
- [157] B. Zong, Q. Song, M. R. Min, W. Cheng, C. Lumezanu, D. Cho, and H. Chen, "Deep autoencoding gaussian mixture model for unsupervised anomaly detection," in *ICLR*, 2018. [7](#) [10](#), [12](#)
- [158] D. Abati, A. Porrello, S. Calderara, and R. Cucchiara, "Latent space autoregression for novelty detection," in *CVPR*, 2019. [7](#) [10](#), [12](#)
- [159] S. Pidhorskyi, R. Almohsen, D. A. Adjeroh, and G. Doretto, "Generative probabilistic novelty detection with adversarial autoencoders," in *NeurIPS*, 2018. [7](#) [10](#), [12](#), [13](#)
- [160] L. Deecke, R. Vandermeulen, L. Ruff, S. Mandt, and M. Kloft, "Image anomaly detection with generative adversarial networks," in *ECML&KDD*, 2018. [7](#) [10](#), [12](#)
- [161] M. Sabokrou, M. Khalooei, M. Fathy, and E. Adeli, "Adversarially learned one-class classifier for novelty detection," in *CVPR*, 2018. [7](#) [10](#), [13](#)
- [162] I. Kobyzhev, S. Prince, and M. Brubaker, "Normalizing flows: An introduction and review of current methods," *TPAMI*, 2020. [7](#) [10](#), [12](#)
- [163] E. Zisselman and A. Tamar, "Deep residual flow for out of distribution detection," in *CVPR*, 2020. [7](#) [10](#)
- [164] D. P. Kingma and P. Dhariwal, "Glow: Generative flow with invertible 1x1 convolutions," *NeurIPS*, 2018. [7](#) [10](#)
- [165] A. Van Oord, N. Kalchbrenner, and K. Kavukcuoglu, "Pixel recurrent neural networks," in *ICML*, 2016. [7](#) [10](#)
- [166] D. Jiang, S. Sun, and Y. Yu, "Revisiting flow generative models for out-of-distribution detection," in *International Conference on Learning Representations*, 2021. [7](#) [10](#)
- [167] E. Nalisnick, A. Matsukawa, Y. W. Teh, D. Gorur, and B. Lakshminarayanan, "Do deep generative models know what they don't know?," in *NeurIPS*, 2018. [7](#) [10](#)
- [168] H. Choi, E. Jang, and A. A. Alemi, "Waic, but why? generative ensembles for robust anomaly detection," *arXiv preprint arXiv:1810.01392*, 2018. [7](#) [10](#)
- [169] P. Kirichenko, P. Izmailov, and A. G. Wilson, "Why normalizing flows fail to detect out-of-distribution data," in *NeurIPS*, 2020. [7](#) [10](#)
- [170] J. Ren, P. J. Liu, E. Fertig, J. Snoek, R. Poplin, M. A. DePristo, J. V. Dillon, and B. Lakshminarayanan, "Likelihood ratios for out-of-distribution detection," in *NeurIPS*, 2019. [7](#) [10](#)
- [171] J. Serrà, D. Álvarez, V. Gómez, O. Slizovskaia, J. F. Núñez, and J. Luque, "Input complexity and out-of-distribution detection with likelihood-based generative models," 2020. [7](#) [10](#)
- [172] Z. Xiao, Q. Yan, and Y. Amit, "Likelihood regret: An out-of-distribution detection score for variational auto-encoder," in *NeurIPS*, 2020. [7](#) [10](#)
- [173] J. Ren, S. Fort, J. Liu, A. G. Roy, S. Padhy, and B. Lakshminarayanan, "A simple fix to mahalanobis distance for improving near-ood detection," *arXiv preprint arXiv:2106.09022*, 2021. [7](#) [10](#)
- [174] E. Techapanurak, M. Suganuma, and T. Okatani, "Hyperparameter-free out-of-distribution detection using cosine similarity," in *ACCV*, 2020. [7](#) [10](#)
- [175] X. Chen, X. Lan, F. Sun, and N. Zheng, "A boundary based out-of-distribution classifier for generalized zero-shot learning," in *ECCV*, 2020. [7](#) [10](#)
- [176] A. Zaeemzadeh, N. Bisagno, Z. Sambugaro, N. Conci, N. Rahnavard, and M. Shah, "Out-of-distribution detection using union of 1-dimensional subspaces," in *CVPR*, 2021. [7](#) [10](#)
- [177] J. Van Amersfoort, L. Smith, Y. W. Teh, and Y. Gal, "Uncertainty estimation using a single deep deterministic neural network," in *ICML*, 2020. [7](#) [10](#)
- [178] H. Huang, Z. Li, L. Wang, S. Chen, B. Dong, and X. Zhou, "Feature space singularity for out-of-distribution detection," *arXiv preprint arXiv:2011.14654*, 2020. [7](#) [10](#)
- [179] Y. Ming, Y. Sun, O. Dia, and Y. Li, "Cider: Exploiting hyperspherical embeddings for out-of-distribution detection," *arXiv preprint arXiv:2203.04450*, 2022. [7](#) [10](#)
- [180] T. Denouden, R. Salay, K. Czarnecki, V. Abdelzad, B. Phan, and S. Vernekar, "Improving reconstruction autoencoder out-of-distribution detection with mahalanobis distance," *arXiv preprint arXiv:1812.02765*, 2018. [7](#) [10](#)
- [181] Y. Zhou, "Rethinking reconstruction autoencoder-based out-of-distribution detection," in *CVPR*, 2022. [7](#) [10](#)
- [182] Y. Yang, R. Gao, and Q. Xu, "Out-of-distribution detection with semantic mismatch under masking," in *ECCV*, 2022. [7](#) [10](#)
- [183] W. J. Scheirer, L. P. Jain, and T. E. Boult, "Probability models for open set recognition," *TPAMI*, 2014. [7](#) [11](#)
- [184] L. P. Jain, W. J. Scheirer, and T. E. Boult, "Multi-class open set recognition using probability of inclusion," in *ECCV*, 2014. [7](#) [11](#)
- [185] A. Bendale and T. E. Boult, "Towards open set deep networks," in *CVPR*, 2016. [7](#) [11](#)
- [186] A. Rozsa, M. Günther, and T. E. Boult, "Adversarial robustness: Softmax versus openmax," in *BMVC*, 2017. [7](#)
- [187] P. Perera and V. M. Patel, "Deep transfer learning for multiple class novelty detection," in *CVPR*, 2019. [7](#) [11](#)
- [188] P. Perera, V. I. Morariu, R. Jain, V. Manjunatha, C. Wigington, V. Ordóñez, and V. M. Patel, "Generative-discriminative feature representations for open-set recognition," in *CVPR*, 2020. [7](#) [11](#)
- [189] X. Sun, H. Ding, C. Zhang, G. Lin, and K.-V. Ling, "M2iosr: Maximal mutual information open set recognition," *arXiv preprint arXiv:2108.02373*, 2021. [7](#) [11](#)
- [190] Z. Ge, S. Demyanov, Z. Chen, and R. Garnavi, "Generative openmax for multi-class open set classification," in *BMVC*, 2017. [7](#) [11](#)
- [191] L. Neal, M. Olson, X. Fern, W.-K. Wong, and F. Li, "Open set learning with counterfactual images," in *ECCV*, 2018. [7](#) [11](#)
- [192] D.-W. Zhou, H.-J. Ye, and D.-C. Zhan, "Learning placeholders for open-set recognition," in *CVPR*, 2021. [7](#) [11](#)
- [193] S. Kong and D. Ramanan, "Opengan: Open-set recognition via open data generation," in *ICCV*, 2021. [7](#) [11](#)
- [194] C. Geng and S. Chen, "Collective decision for open set recognition," *TKDE*, 2020. [7](#) [11](#)
- [195] J. Jang and C. O. Kim, "One-vs-rest network-based deep probability model for open set recognition," *arXiv preprint arXiv:2004.08067*, 2020. [7](#) [11](#)
- [196] P. Schlachter, Y. Liao, and B. Yang, "Open-set recognition using intra-class splitting," in *EUSIPCO*, 2019. [7](#) [11](#)
- [197] M. Masana, I. Ruiz, J. Serrat, J. van de Weijer, and A. M. Lopez, "Metric learning for novelty and anomaly detection," in *BMVC*, 2018. [7](#) [11](#)
- [198] Y. Shu, Y. Shi, Y. Wang, T. Huang, and Y. Tian, "p-odn: prototype-based open deep network for open set recognition," *Scientific reports*, 2020. [7](#) [11](#)

- [199] B. Liu, H. Kang, H. Li, G. Hua, and N. Vasconcelos, "Few-shot open-set recognition using meta-learning," in *CVPR*, 2020. [7, 11](#)
- [200] G. Chen, L. Qiao, Y. Shi, P. Peng, J. Li, T. Huang, S. Pu, and Y. Tian, "Learning open set network with discriminative reciprocal points," in *ECCV*, 2020. [7, 11](#)
- [201] R. Yoshihashi, W. Shao, R. Kawakami, S. You, M. Iida, and T. Naemura, "Classification-reconstruction learning for open-set recognition," in *CVPR*, 2019. [7, 11](#)
- [202] A. Cao, Y. Luo, and D. Klabjan, "Open-set recognition with gaussian mixture variational autoencoders," *AAAI*, 2020. [7, 11](#)
- [203] P. R. M. Júnior, R. M. De Souza, R. d. O. Werneck, B. V. Stein, D. V. Pazinato, W. R. de Almeida, O. A. Penatti, R. d. S. Torres, and A. Rocha, "Nearest neighbors distance ratio open-set classifier," *Machine Learning*, 2017. [7, 11](#)
- [204] H. Zhang and V. M. Patel, "Sparse representation-based open set recognition," *TPAMI*, 2016. [7, 11](#)
- [205] P. Bodesheim, A. Freytag, E. Rodner, M. Kemmler, and J. Denzler, "Kernel null space methods for novelty detection," in *CVPR*, 2013. [7, 11](#)
- [206] J. Liu, Z. Lian, Y. Wang, and J. Xiao, "Incremental kernel null space discriminant analysis for novelty detection," in *CVPR*, 2017. [7, 11](#)
- [207] P. Oza and V. M. Patel, "C2ae: Class conditioned auto-encoder for open-set recognition," in *CVPR*, 2019. [7, 11](#)
- [208] X. Sun, Z. Yang, C. Zhang, K.-V. Ling, and G. Peng, "Conditional gaussian distribution learning for open set recognition," in *CVPR*, 2020. [7, 11](#)
- [209] Z. Yue, T. Wang, Q. Sun, X.-S. Hua, and H. Zhang, "Counterfactual zero-shot and open-set visual recognition," in *CVPR*, 2021. [7, 11](#)
- [210] R. Shao, P. Perera, P. C. Yuen, and V. M. Patel, "Open-set adversarial defense," in *ECCV*, 2020. [7, 11](#)
- [211] G. Danuser and M. Stricker, "Parametric model fitting: From inlier characterization to outlier detection," *TPAMI*, 1998. [7, 12](#)
- [212] C. Leys, O. Klein, Y. Dominicy, and C. Ley, "Detecting multivariate outliers: Use a robust variant of the mahalanobis distance," *Journal of Experimental Social Psychology*, 2018. [7, 12](#)
- [213] E. Eskin, "Anomaly detection over noisy data using learned probability distributions," in *ICML*, 2000. [7, 12](#)
- [214] M. Turcotte, J. Moore, N. Heard, and A. McPhall, "Poisson factorization for peer-based anomaly detection," in *IEEE Conference on Intelligence and Security Informatics (ISI)*, 2016. [7, 12](#)
- [215] A. J. Izenman, "Review papers: Recent developments in non-parametric density estimation," *Journal of the American Statistical Association*, 1991. [7, 12](#)
- [216] W. Hu, J. Gao, B. Li, O. Wu, J. Du, and S. Maybank, "Anomaly detection using local kernel density estimation and context-based regression," *TKDE*, 2018. [7, 12](#)
- [217] D. Rezende and S. Mohamed, "Variational inference with normalizing flows," in *ICML*, 2015. [7, 12](#)
- [218] C.-L. Li, K. Sohn, J. Yoon, and T. Pfister, "Cutpaste: Self-supervised learning for anomaly detection and localization," in *CVPR*, 2021. [7, 12](#)
- [219] M. Salehi, N. Sadjadi, S. Baselizadeh, M. H. Rohban, and H. R. Rabiee, "Multiresolution knowledge distillation for anomaly detection," in *CVPR*, 2021. [7, 12](#)
- [220] A. K. Menon and R. C. Williamson, "A loss framework for calibrated anomaly detection," in *NeurIPS*, 2018. [7, 12](#)
- [221] H. Zhao, W. Zhou, D. Chen, T. Wei, W. Zhang, and N. Yu, "Multi-attentional deepfake detection," in *CVPR*, 2021. [7, 12](#)
- [222] P. Yang, J. Chen, C.-J. Hsieh, J.-L. Wang, and M. Jordan, "Ml-loo: Detecting adversarial examples with feature attribution," in *AAAI*, 2020. [7, 12](#)
- [223] M. Du, S. Pentylala, Y. Li, and X. Hu, "Towards generalizable deepfake detection with locality-aware autoencoder," in *CIKM*, 2020. [7, 12](#)
- [224] S. Zhai, Y. Cheng, W. Lu, and Z. Zhang, "Deep structured energy based models for anomaly detection," in *ICML*, 2016. [7, 12](#)
- [225] A. Hyvärinen and P. Dayan, "Estimation of non-normalized statistical models by score matching," 2005. [7, 12](#)
- [226] M. Welling and Y. W. Teh, "Bayesian learning via stochastic gradient langevin dynamics," in *ICML*, 2011. [7, 12](#)
- [227] H. Wang, X. Wu, Z. Huang, and E. P. Xing, "High-frequency component helps explain the generalization of convolutional neural networks," in *CVPR*, 2020. [7, 12](#)
- [228] A. Ilyas, S. Santurkar, D. Tsipras, L. Engstrom, B. Tran, and A. Madry, "Adversarial examples are not bugs, they are features," in *NeurIPS*, 2019. [7, 12](#)
- [229] G. Chen, P. Peng, L. Ma, J. Li, L. Du, and Y. Tian, "Amplitude-phase recombination: Rethinking robustness of convolutional neural networks in frequency domain," *ICCV*, 2021. [7, 12](#)
- [230] H. Liu, X. Li, W. Zhou, Y. Chen, Y. He, H. Xue, W. Zhang, and N. Yu, "Spatial-phase shallow learning: rethinking face forgery detection in frequency domain," in *CVPR*, 2021. [7, 13](#)
- [231] A. Adler, M. Elad, Y. Hel-Or, and E. Rivlin, "Sparse coding with anomaly detection," *Journal of Signal Processing Systems*, 2015. [7, 13](#)
- [232] A. Li, Z. Miao, Y. Cen, and Y. Cen, "Anomaly detection using sparse reconstruction in crowded scenes," *Multimedia Tools and Applications*, 2017. [7, 13](#)
- [233] X. Mo, V. Monga, R. Bala, and Z. Fan, "Adaptive sparse representations for video anomaly detection," *IEEE Transactions on Circuits and Systems for Video Technology*, 2013. [7, 13](#)
- [234] Y. Xiao, H. Wang, W. Xu, and J. Zhou, "L1 norm based kPCA for novelty detection," *Pattern Recognition*, 2013. [7, 13](#)
- [235] K. Jiang, W. Xie, J. Lei, T. Jiang, and Y. Li, "Lren: Low-rank embedded network for sample-free hyperspectral anomaly detection," in *AAAI*, 2021. [7, 13](#)
- [236] Z. Chen, C. K. Yeo, B. S. Lee, and C. T. Lau, "Autoencoder-based network anomaly detection," in *Wireless Telecommunications Symposium*, 2018. [7, 13](#)
- [237] J. An and S. Cho, "Variational autoencoder based anomaly detection using reconstruction probability," *Special Lecture on IE*, 2015. [7, 13](#)
- [238] H. Zenati, C. S. Foo, B. Lecouat, G. Manek, and V. R. Chandrasekhar, "Efficient gan-based anomaly detection," in *ICLR-W*, 2018. [7, 13](#)
- [239] W. Liu, W. Luo, D. Lian, and S. Gao, "Future frame prediction for anomaly detection-a new baseline," in *CVPR*, 2018. [7, 13](#)
- [240] D. Gong, L. Liu, V. Le, B. Saha, M. R. Mansour, S. Venkatesh, and A. v. d. Hengel, "Memorizing normality to detect anomalies: Memory-augmented deep autoencoder for unsupervised anomaly detection," in *CVPR*, 2019. [7, 13](#)
- [241] H. Park, J. Noh, and B. Ham, "Learning memory-guided normality for anomaly detection," in *CVPR*, 2020. [7, 13](#)
- [242] C.-H. Lai, D. Zou, and G. Lerman, "Robust subspace recovery layer for unsupervised anomaly detection," *ICLR*, 2020. [7, 13](#)
- [243] X. Yan, H. Zhang, X. Xu, X. Hu, and P.-A. Heng, "Learning semantic context from normal samples for unsupervised anomaly detection," in *AAAI*, 2021. [7, 13](#)
- [244] D. T. Nguyen, Z. Lou, M. Klar, and T. Brox, "Anomaly detection with multiple-hypotheses predictions," in *ICML*, 2019. [7, 13](#)
- [245] K. Tian, S. Zhou, J. Fan, and J. Guan, "Learning competitive and discriminative reconstructions for anomaly detection," in *AAAI*, 2019. [7, 13](#)
- [246] K. Zhou, Y. Xiao, J. Yang, J. Cheng, W. Liu, W. Luo, Z. Gu, J. Liu, and S. Gao, "Encoding structure-texture relation with p-net for anomaly detection in retinal images," in *ECCV*, 2020. [7, 13](#)
- [247] X. Han, X. Chen, and L.-P. Liu, "Gan ensemble for anomaly detection," *arXiv preprint arXiv:2012.07988*, 2020. [7, 13](#)
- [248] G. Kwon, M. Prabhushankar, D. Temel, and G. AlRegib, "Back-propagated gradient representations for anomaly detection," in *ECCV*, 2020. [7, 13](#)
- [249] J. Tian, M. H. Azarian, and M. Pecht, "Anomaly detection using self-organizing maps-based k-nearest neighbor algorithm," in *PHM Society European Conference*, 2014. [7, 13](#)
- [250] G. Münz, S. Li, and G. Carle, "Traffic anomaly detection using k-means clustering," in *GI/ITG Workshop MMBnet*, 2007. [7, 13](#)
- [251] I. Syarif, A. Prugel-Bennett, and G. Wills, "Unsupervised clustering approach for network anomaly detection," in *International conference on networked digital technologies*, 2012. [7, 13](#)
- [252] D. M. J. Tax, "One-class classification: Concept learning in the absence of counter-examples," 2002. [7, 13, 14](#)
- [253] L. Ruff, R. Vandermeulen, N. Goernitz, L. Deecke, S. A. Siddiqui, A. Binder, E. Müller, and M. Kloft, "Deep one-class classification," in *ICML*, 2018. [7, 13, 14](#)
- [254] T. Reiss, N. Cohen, L. Bergman, and Y. Hoshen, "Panda: Adapting pretrained features for anomaly detection and segmentation," in *CVPR*, 2021. [7, 13](#)
- [255] J. Wang and A. Cherian, "Gods: Generalized one-class discriminative subspaces for anomaly detection," in *CVPR*, 2019. [7, 13](#)
- [256] B. Zhang and W. Zuo, "Learning from positive and unlabeled examples: A survey," in *International Symposiums on Information Processing*, 2008. [7, 13](#)

- [257] J. Bekker and J. Davis, "Learning from positive and unlabeled data: A survey," *Machine Learning*, 2020. **7**, 13
- [258] K. Jaskie and A. Spanias, "Positive and unlabeled learning algorithms and applications: A survey," in *International Conference on Information, Intelligence, Systems and Applications*, 2019. **7**, 13
- [259] C. Wang, C. Ding, R. F. Meraz, and S. R. Holbrook, "Psol: a positive sample only learning algorithm for finding non-coding rna genes," *Bioinformatics*, 2006. **7**, 13
- [260] X. Li and B. Liu, "Learning to classify texts using positive and unlabeled data," in *IJCAI*, 2003. **7**, 13
- [261] B. Zhang and W. Zuo, "Reliable negative extracting based on knn for learning from positive and unlabeled examples," *Journal of Computers*, 2009. **7**, 13
- [262] S. Chaudhari and S. Shevade, "Learning from positive and unlabelled examples using maximum margin clustering," in *ICONIP*, 2012. **7**, 13
- [263] L. Liu and T. Peng, "Clustering-based method for positive and unlabeled text categorization enhanced by improved tfidf," *Journal of Information Science and Engineering*, 2014. **7**, 13
- [264] F. He, T. Liu, G. I. Webb, and D. Tao, "Instance-dependent pu learning by bayesian optimal relabeling," *arXiv preprint arXiv:1808.02180*, 2018. **7**, 13
- [265] A. Menon, B. Van Rooyen, C. S. Ong, and B. Williamson, "Learning from corrupted binary labels via class-probability estimation," in *ICML*, 2015. **7**, 13
- [266] C. Scott, "A rate of convergence for mixture proportion estimation, with application to learning from noisy labels," in *Artificial Intelligence and Statistics*, 2015. **7**, 13
- [267] J.-X. Zhong, N. Li, W. Kong, S. Liu, T. H. Li, and G. Li, "Graph convolutional label noise cleaner: Train a plug-and-play action classifier for anomaly detection," in *CVPR*, 2019. **7**, 13
- [268] K. Tian, S. Zhou, J. Fan, and J. Guan, "Learning competitive and discriminative reconstructions for anomaly detection," in *AAAI*, 2019. **7**, 13
- [269] F. T. Liu, K. M. Ting, and Z.-H. Zhou, "Isolation forest," in *ICDM*, 2008. **7**, 14
- [270] J. Tack, S. Mo, J. Jeong, and J. Shin, "Csi: Novelty detection via contrastive learning on distributionally shifted instances," in *NeurIPS*, 2020. **7**, 14
- [271] L. Bergman and Y. Hoshen, "Classification-based anomaly detection for general data," in *ICLR*, 2020. **7**, 14
- [272] I. Golani and R. El-Yaniv, "Deep anomaly detection using geometric transformations," in *NeurIPS*, 2018. **7**, 14
- [273] M.-I. Georgescu, A. Barbalau, R. T. Ionescu, F. S. Khan, M. Popescu, and M. Shah, "Anomaly detection in video via self-supervised and multi-task learning," in *CVPR*, 2021. **7**, 14
- [274] X. Wan, W. Wang, J. Liu, and T. Tong, "Estimating the sample mean and standard deviation from the sample size, median, range and/or interquartile range," *BMC medical research methodology*, 2014. **7**, 14
- [275] M. M. Breunig, H.-P. Kriegel, R. T. Ng, and J. Sander, "Lof: identifying density-based local outliers," in *SIGMOD*, 2000. **7**, 14
- [276] J. Tang, Z. Chen, A. W.-C. Fu, and D. W. Cheung, "Enhancing effectiveness of outlier detections for low density patterns," in *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, 2002. **7**, 14
- [277] H.-P. Kriegel, P. Kröger, E. Schubert, and A. Zimek, "Loop: local outlier probabilities," in *Proceedings of the 18th ACM conference on Information and knowledge management*, 2009. **7**, 14
- [278] E. Schubert, A. Zimek, and H.-P. Kriegel, "Local outlier detection reconsidered: a generalized view on locality with applications to spatial, video, and network outlier detection," *Data mining and knowledge discovery*, 2014. **7**, 14
- [279] M. A. Fischler and R. C. Bolles, "Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography," *Communications of the ACM*, 1981. **7**, 14
- [280] P. J. Rousseeuw and M. Hubert, "Robust statistics for outlier detection," *Wiley interdisciplinary reviews: Data mining and knowledge discovery*, 2011. **7**, 14
- [281] V. Sharan, P. Gopalan, and U. Wieder, "Efficient anomaly detection via matrix sketching," *NeurIPS*, 2018. **7**, 14
- [282] M. Ester, H.-P. Kriegel, J. Sander, X. Xu, et al., "A density-based algorithm for discovering clusters in large spatial databases with noise," in *KDD*, 1996. **7**, 14
- [283] U. Rebbapragada and C. E. Brodley, "Class noise mitigation through instance weighting," in *ECML*, 2007. **7**, 14
- [284] L. Akoglu, H. Tong, and D. Koutra, "Graph based anomaly detection and description: a survey," *Data mining and knowledge discovery*, 2015. **7**, 14
- [285] C. C. Noble and D. J. Cook, "Graph-based anomaly detection," in *SICKDD*, 2003. **7**, 14
- [286] Y. Kou, C.-T. Lu, and R. F. Dos Santos, "Spatial outlier detection: a graph-based approach," in *19th IEEE International Conference on Tools with Artificial Intelligence (ICTAI)*, 2007. **7**, 14
- [287] Z. Mingqiang, H. Hui, and W. Qian, "A graph-based clustering algorithm for anomaly intrusion detection," in *International Conference on Computer Science & Education (ICCSE)*, 2012. **7**, 14
- [288] J. Yang, W. Chen, L. Feng, X. Yan, H. Zheng, and W. Zhang, "Webly supervised image classification with metadata: Automatic noisy label correction via visual-semantic graph," in *ACM Multimedia*, 2020. **7**, 14
- [289] Y. Li, J. Yang, Y. Song, L. Cao, J. Luo, and L.-J. Li, "Learning from noisy labels with distillation," in *CVPR*, 2017. **7**, 14
- [290] D. T. Nguyen, C. K. Mummadipati, T. P. N. Ngo, T. H. P. Nguyen, L. Beguel, and T. Brox, "Self: Learning to filter noisy labels with self-ensembling," in *ICLR*, 2020. **7**, 14
- [291] B. Han, Q. Yao, X. Yu, G. Niu, M. Xu, W. Hu, I. Tsang, and M. Sugiyama, "Co-teaching: Robust training of deep neural networks with extremely noisy labels," in *NIPS*, 2018. **7**, 14
- [292] J. Yang, L. Feng, W. Chen, X. Yan, H. Zheng, P. Luo, and W. Zhang, "Webly supervised image classification with self-contained confidence," in *ECCV*, 2020. **7**, 14
- [293] P. Morteza and Y. Li, "Provable guarantees for understanding out-of-distribution detection," in *AAAI*, 2022. **8**, 10, 16
- [294] K. Bibas, M. Feder, and T. Hassner, "Single layer predictive normalized maximum likelihood for out-of-distribution detection," *NeurIPS*, 2021. **8**
- [295] X. Dong, J. Guo, W.-T. T. Ang Li23, C. Liu, and H. Kung, "Neural mean discrepancy for efficient out-of-distribution detection," in *CVPR*, 2022. **8**
- [296] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al., "Learning transferable visual models from natural language supervision," in *ICML*, 2021. **9**
- [297] E. T. Jaynes, "Bayesian methods: General background," 1986. **9**
- [298] R. M. Neal, *Bayesian learning for neural networks*. 2012. **9**
- [299] D. Gamerman and H. F. Lopes, *Markov chain Monte Carlo: stochastic simulation for Bayesian inference*. CRC Press, 2006. **9**
- [300] D. J. C. Mackay, *Bayesian methods for adaptive models*. PhD thesis, California Institute of Technology, 1992. **9**
- [301] A. Y. Foong, Y. Li, J. M. Hernández-Lobato, and R. E. Turner, "'in-between' uncertainty in bayesian neural networks," in *ICML-W*, 2020. **9**
- [302] C. Peterson and E. Hartman, "Explorations of the mean field theory learning algorithm," *Neural Networks*, 1989. **9**
- [303] F. Wenzel, K. Roth, B. S. Veeling, J. Świątkowski, L. Tran, S. Mandt, J. Snoek, T. Salimans, R. Jenatton, and S. Nowozin, "How good is the bayes posterior in deep neural networks really?", in *ICML*, 2020. **9**
- [304] A. Gelman, "Objections to bayesian statistics," *Bayesian Analysis*, 2008. **9**
- [305] T. G. Dietterich, "Ensemble methods in machine learning," in *International workshop on multiple classifier systems*, 2000. **9**
- [306] W. J. Maddox, P. Izmailov, T. Garipov, D. P. Vetrov, and A. G. Wilson, "A simple baseline for bayesian uncertainty in deep learning," *Advances in Neural Information Processing Systems*, vol. 32, pp. 13153–13164, 2019. **9**
- [307] K. Kim, J. Shin, and H. Kim, "Locally most powerful bayesian test for out-of-distribution detection using deep generative models," *NeurIPS*, 2021. **9**
- [308] J. Yang, K. Zhou, and Z. Liu, "Full-spectrum out-of-distribution detection," *arXiv preprint arXiv:2204.05306*, 2022. **10**, 15
- [309] E. D. C. Gomes, F. Alberge, P. Duhamel, and P. Piantanida, "Igeood: An information geometry approach to out-of-distribution detection," in *ICLR*, 2022. **10**
- [310] H. Wang, Z. Li, L. Feng, and W. Zhang, "Vim: Out-of-distribution with virtual-logit matching," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. **10**
- [311] R. L. Smith, "Extreme value theory," *Handbook of applicable mathematics*, 1990. **11**
- [312] E. Castillo, *Extreme value theory in engineering*. Elsevier, 2012. **11**
- [313] H. Zhang, A. Li, J. Guo, and Y. Guo, "Hybrid models for open set recognition," in *ECCV*, 2020. **11**

- [314] S. Vaze, K. Han, A. Vedaldi, and A. Zisserman, "Open-set recognition: A good closed-set classifier is all you need," in *ICLR*, 2022. [11](#), [15](#)
- [315] R. De Maesschalck, D. Jouan-Rimbaud, and D. L. Massart, "The mahalanobis distance," *Chemometrics and intelligent laboratory systems*, 2000. [12](#), [14](#)
- [316] R. A. Redner and H. F. Walker, "Mixture densities, maximum likelihood and the em algorithm," *SIAM review*, 1984. [12](#)
- [317] J. Van Ryzin, "A histogram method of density estimation," *Communications in Statistics-Theory and Methods*, 1973. [12](#)
- [318] M. Xie, J. Hu, and B. Tian, "Histogram-based online anomaly detection in hierarchical wireless sensor networks," in *ICTSPCC*, 2012. [12](#)
- [319] A. Kind, M. P. Stoecklin, and X. Dimitropoulos, "Histogram-based traffic anomaly detection," *IEEE Transactions on Network and Service Management*, 2009. [12](#)
- [320] M. Goldstein and A. Dengel, "Histogram-based outlier score (hbos): A fast unsupervised anomaly detection algorithm," *KI-2012: Poster and Demo Track*, 2012. [12](#)
- [321] E. Parzen, "On estimation of a probability density function and mode," *The annals of mathematical statistics*, 1962. [12](#)
- [322] M. Desforges, P. Jacob, and J. Cooper, "Applications of probability density estimation to the detection of abnormal conditions in engineering," *Proceedings of the institution of mechanical engineers*, 1998. [12](#)
- [323] K. Nishiya, J. Hasegawa, and T. Koike, "Dynamic state estimation including anomaly detection and identification for power systems," in *IEE proceedings C (generation, transmission and distribution)*, 1982. [12](#)
- [324] P. Helman and G. Liepins, "Statistical foundations of audit trail analysis for the detection of computer misuse," *IEEE Transactions on software engineering*, 1993. [12](#)
- [325] P. D. Talagala, R. J. Hyndman, and K. Smith-Miles, "Anomaly detection in high-dimensional data," *Journal of Computational and Graphical Statistics*, 2020. [12](#)
- [326] A. Patcha and J.-M. Park, "An overview of anomaly detection techniques: Existing solutions and latest technological trends," *Computer networks*, 2007. [12](#)
- [327] S. Wu, B. E. Moore, and M. Shah, "Chaotic invariants of lagrangian particle trajectories for anomaly detection in crowded scenes," in *CVPR*, 2010. [12](#)
- [328] M. A. Kramer, "Nonlinear principal component analysis using autoassociative neural networks," *AIChE journal*, 1991. [12](#)
- [329] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013. [12](#)
- [330] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," *NIPS*, 2014. [12](#)
- [331] J. Ngiam, Z. Chen, P. W. Koh, and A. Y. Ng, "Learning deep energy models," in *ICML*, 2011. [12](#)
- [332] J. Vanschoren, "Meta-learning: A survey," *arXiv preprint arXiv:1810.03548*, 2018. [13](#)
- [333] T. Hospedales, A. Antoniou, P. Micaelli, and A. Storkey, "Meta-learning in neural networks: A survey," *arXiv preprint arXiv:2004.05439*, 2020. [13](#)
- [334] D. Wettschereck, "A study of distance-based machine learning algorithms," 1994. [13](#)
- [335] N. Görnitz, M. Kloft, K. Rieck, and U. Brefeld, "Toward supervised anomaly detection," *Journal of Artificial Intelligence Research*, 2013. [13](#)
- [336] S. S. Khan and M. G. Madden, "A survey of recent trends in one class classification," in *Irish conference on artificial intelligence and cognitive science*, 2009. [13](#)
- [337] D. G. Altman and J. M. Bland, "Standard deviations and standard errors," *BMJ*, 2005. [14](#)
- [338] C. Ley, C. Ley, O. Klein, P. Bernard, and L. Licata, "Detecting outliers: Do not use standard deviation around the mean, use absolute deviation around the median," *Journal of experimental social psychology*, 2013. [14](#)
- [339] X. Yang, L. J. Latecki, and D. Pokrajac, "Outlier detection with globally optimal exemplar-based gmm," in *SIAM*, 2009. [14](#)
- [340] M. Sugiyama and K. Borgwardt, "Rapid distance-based outlier detection via sampling," *NIPS*, 2013. [14](#)
- [341] G. H. Orair, C. H. Teixeira, W. Meira Jr, Y. Wang, and S. Parthasarathy, "Distance-based outlier detection: consolidation and renewed bearing," *Proceedings of the VLDB Endowment*, 2010. [14](#)
- [342] V. Hautamaki, I. Karkkainen, and P. Franti, "Outlier detection using k-nearest neighbour graph," in *ICPR*, 2004. [14](#)
- [343] F. Muhlenbach, S. Lallich, and D. A. Zighed, "Identifying and handling mislabelled instances," *Journal of Intelligent Information Systems*, 2004. [14](#)
- [344] W. Liu, J. He, and S.-F. Chang, "Large graph construction for scalable semi-supervised learning," in *ICML*, 2010. [14](#)
- [345] Z.-F. Wu, T. Wei, J. Jiang, C. Mao, M. Tang, and Y.-F. Li, "Ngc: A unified framework for learning with open-world noisy data," in *ICCV*, 2021. [14](#), [16](#)
- [346] Y. Ming, H. Yin, and Y. Li, "On the impact of spurious correlation for out-of-distribution detection," in *AAAI*, 2022. [15](#)
- [347] P. Panareda Busto and J. Gall, "Open set domain adaptation," in *ICCV*, 2017. [15](#)
- [348] Y. Shu, Z. Cao, C. Wang, J. Wang, and M. Long, "Open domain generalization with domain-augmented meta-learning," in *CVPR*, 2021. [15](#)
- [349] J. Li, C. Xiong, and S. C. Hoi, "Mopro: Webly supervised learning with momentum prototypes," *ICLR*, 2021. [15](#)
- [350] R. Kaur, K. Sridhar, S. Park, S. Jha, A. Roy, O. Sokolsky, and I. Lee, "Codit: Conformal out-of-distribution detection in time-series data," *arXiv e-prints*, 2022. [16](#)
- [351] V. D. Nguyen, "Out-of-distribution detection for lidar-based 3d object detection," Master's thesis, University of Waterloo, 2022. [16](#)
- [352] S. Esmaeilpour, B. Liu, E. Robertson, and L. Shu, "Zero-shot out-of-distribution detection based on the pretrained model clip," in *AAAI*, 2022. [16](#)
- [353] G. Shalev, G.-L. Shalev, and J. Keshet, "A baseline for detecting out-of-distribution examples in image captioning," *arXiv preprint arXiv:2207.05418*, 2022. [16](#)