

---

# Is Out-of-distribution Detection Learnable?

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       Supervised learning aims to train a classifier under the assumption that training and  
2       test data are from the same distribution. To ease the above assumption, researchers  
3       have studied a more realistic setting: *out-of-distribution* (OOD) detection, where  
4       test data may come from classes that are unknown during training (*i.e.*, OOD data).  
5       Due to the unavailability and diversity of OOD data, good generalization ability  
6       is crucial for effective OOD detection algorithms. To study the generalization of  
7       OOD detection, in this paper, we investigate the *probably approximately correct*  
8       (PAC) learning theory of OOD detection, which is proposed by researchers as an  
9       *open problem*. First, we find a necessary condition for the learnability of OOD  
10       detection. Then, using this condition, we prove several impossibility theorems for  
11       the learnability of OOD detection under some scenarios. Although the impossibil-  
12       ity theorems are frustrating, we find that some conditions of these impossibility  
13       theorems may not hold in some practical scenarios. Based on this observation, we  
14       next give several necessary and sufficient conditions to characterize the learnability  
15       of OOD detection in some practical scenarios. Lastly, we also offer theoretical  
16       supports for several representative OOD detection works based on our OOD theory.

## 17   1 Introduction

18   The success of supervised learning is established on an implicit assumption that training and test data  
19   share a same distribution, *i.e.*, *in-distribution* (ID) [1, 2, 3, 4]. However, test data distribution in many  
20   real-world scenarios may violate the assumption and, instead, contain *out-of-distribution* (OOD) data  
21   whose labels have not been seen during the training process [5, 6]. To mitigate the risk of OOD data,  
22   researchers have considered a more practical learning scenario: OOD detection which determines  
23   whether an input is ID/OOD, while classifying the ID data into respective classes. OOD detection has  
24   shown great potential to ensure the reliable deployment of machine learning models in the real world.  
25   A rich line of algorithms have been developed to empirically address the OOD detection problem  
26   [6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20]. However, very few works study theory of OOD  
27   detection, which hinders the rigorous path forward for the field. This paper aims to bridge the gap.

28   In this paper, we provide a theoretical framework to understand the learnability of the OOD detection  
29   problem. We investigate the probably approximately correct (PAC) learning theory of OOD detection,  
30   which is posed as an open problem to date. Unlike the classical PAC learning theory in a supervised  
31   setting, our problem setting is fundamentally challenging due to the *absence of OOD data* in training.  
32   In many real-world scenarios, OOD data can be diverse and priori-unknown. Given this, we study  
33   whether there exists an algorithm that can be used to detect various OOD data instead of merely some  
34   specified OOD data. Such is the significance of studying the learning theory for OOD detection [4].  
35   This motivates our question: *is OOD detection agnostic PAC learnable? i.e., is there the agnostic*  
36   *PAC learning theory to guarantee the generalization ability of OOD detection?*

37   To investigate the learning theory, we mainly focus on two basic spaces: domain space and hypothesis  
38   space. The domain space is a space consisting of some distributions, and the hypothesis space is a

39 space consisting of some classifiers. Existing agnostic PAC theories in supervised learning [21, 22]  
 40 are distribution-free, *i.e.*, the domain space consists of all domains. Yet, in Theorem 4, we shows that  
 41 the learning theory of OOD detection is not distribution-free. In fact, we discover that OOD detection  
 42 is learnable only if the domain space and the hypothesis space satisfy some special conditions, *e.g.*,  
 43 Conditions 1 and 3. Notably, there are many conditions and theorems in existing learning theories  
 44 and many OOD detection algorithms in the literature. Thus, it is very difficult to analyze the relation  
 45 between these theories and algorithms, and explore useful conditions to ensure the learnability of  
 46 OOD detection, especially when we have to explore them *from the scratch*. Thus, the main aim of our  
 47 paper is to study these essential conditions. From these essential conditions, we can know *when* OOD  
 48 detection can be successful in practical scenarios. We restate our question and goal in following:

49 *Given hypothesis spaces and several representative domain spaces, what are  
 the conditions to ensure the learnability of OOD detection? If possible, we  
 hope that these conditions are necessary and sufficient in some scenarios.*

50 **Main Results.** We investigate the learnability of OOD detection starting from the largest space—the  
 51 total space, and give a necessary condition (Condition 1) for the learnability. However, we find that  
 52 the overlap between ID and OOD data may result in that the necessary condition does not hold.  
 53 Therefore, we give an impossibility theorem to demonstrate that OOD detection fails in the total  
 54 space (Theorem 4). Next, we study OOD detection in the separate space, where there are no overlaps  
 55 between the ID and OOD data. Unfortunately, there still exists impossibility theorem (Theorem 5),  
 56 which demonstrates that OOD detection is not learnable in the separate space under some conditions.

57 Although the impossibility theorems obtained in the separate space are frustrating, we find that some  
 58 conditions of these impossibility theorems may not hold in some practical scenarios. Based on this  
 59 observation, we give several necessary and sufficient conditions to characterize the learnability of  
 60 OOD detection in the separate space (Theorems 6 and 10). Especially, when our model is based on  
 61 *fully-connected neural network* (FCNN), OOD detection is learnable in the separate space if and  
 62 only if the feature space is finite. Furthermore, we investigate the learnability of OOD detection in  
 63 other more practical domain spaces, *e.g.*, the finite-ID-distribution space (Theorem 8) and the density-  
 64 based space (Theorem 9). By studying the finite-ID-distribution space, we discover a compatibility  
 65 condition (Condition 3) that is a necessary and sufficient condition for this space. Next, we further  
 66 investigate the compatibility condition in the density-based space, and find that such condition is also  
 67 the necessary and sufficient condition in some practical scenarios (Theorem 11).

68 **Implications and Impacts of Theory.** Our study is not of purely theoretical interest; it has also  
 69 practical impacts. First, when we design OOD detection algorithms, we normally only have finite  
 70 ID datasets, corresponding to the finite-ID-distribution space. In this case, Theorem 8 provides the  
 71 necessary and sufficient condition to the success of OOD detection. Second, our theory also provides  
 72 theoretical support (Theorems 10 and 11) for several representative OOD detection works [7, 8, 23].  
 73 Third, our theory shows that OOD detection can be addressed in image-based distributions as long as  
 74 ID images have clearly different semantic meanings from OOD images. Fourth, we should not expect  
 75 a universally working algorithm. It is necessary to design different algorithms in different scenarios.

## 76 2 Learning Setups

77 We start by introducing the necessary concepts and notations for our theoretical framework. Given  
 78 a feature space  $\mathcal{X} \subset \mathbb{R}^d$  and a label space  $\mathcal{Y} := \{1, \dots, K\}$ , we have an ID joint distribution  $D_{X_I Y_I}$   
 79 over  $\mathcal{X} \times \mathcal{Y}$ , where  $X_I \in \mathcal{X}$  and  $Y_I \in \mathcal{Y}$  are random variables. We also have an OOD joint  
 80 distribution  $D_{X_O Y_O}$ , where  $X_O$  is a random variable from  $\mathcal{X}$ , but  $Y_O$  is a random variable whose  
 81 outputs do not belong to  $\mathcal{Y}$ . During testing, we will meet a mixture of ID and OOD joint distributions:  
 82  $D_{XY} := (1 - \pi^{\text{out}})D_{X_I Y_I} + \pi^{\text{out}}D_{X_O Y_O}$ , and can only observe the marginal distribution  $D_X :=$   
 83  $(1 - \pi^{\text{out}})D_{X_I} + \pi^{\text{out}}D_{X_O}$ , where the constant  $\pi^{\text{out}} \in [0, 1)$  is an unknown class-prior probability.

84 **Problem 1** (OOD Detection [4]). *Given an ID joint distribution  $D_{X_I Y_I}$  and a training data  $S :=$   
 85  $\{(\mathbf{x}^1, y^1), \dots, (\mathbf{x}^n, y^n)\}$  drawn independent and identically distributed from  $D_{X_I Y_I}$ , the aim of OOD  
 86 detection is to train a classifier  $f$  by using the training data  $S$  such that, for any test data  $\mathbf{x}$  drawn  
 87 from the mixed marginal distribution  $D_X$ : 1) if  $\mathbf{x}$  is an observation from  $D_{X_I}$ ,  $f$  can classify  $\mathbf{x}$  into  
 88 correct ID classes; and 2) if  $\mathbf{x}$  is an observation from  $D_{X_O}$ ,  $f$  can detect  $\mathbf{x}$  as OOD data.*

89 According to the survey [4], when  $K > 1$ , OOD detection is also known as the open-set recognition  
 90 or open-set learning [24, 25]; and when  $K = 1$ , OOD detection reduces to one-class novelty detection  
 91 and semantic anomaly detection [26, 27, 28].

92 **OOD Label and Domain Space.** Based on Problem 1, we know it is not necessary to classify OOD  
 93 data into the correct OOD classes. Without loss of generality, let all OOD data be allocated to one big  
 94 OOD class, *i.e.*,  $Y_O = K + 1$  [24, 29]. To investigate the agnostic PAC learnability of OOD detection,  
 95 we define a domain space  $\mathcal{D}_{XY}$ , which is a set consisting of some joint distributions  $D_{XY}$  mixed by  
 96 some ID joint distributions and some OOD joint distributions. In this paper, the joint distribution  
 97  $D_{XY}$  mixed by ID joint distribution  $D_{X_I Y_I}$  and OOD joint distribution  $D_{X_O Y_O}$  is called **domain**.

98 **Hypothesis Spaces and Scoring Function Spaces.** A hypothesis space  $\mathcal{H}$  is a subset of function  
 99 space, *i.e.*,  $\mathcal{H} \subset \{h : \mathcal{X} \rightarrow \mathcal{Y} \cup \{K + 1\}\}$ . We set  $\mathcal{H}^{\text{in}} \subset \{h : \mathcal{X} \rightarrow \mathcal{Y}\}$  to the ID hypothesis space.  
 100 We also define  $\mathcal{H}^{\text{b}} \subset \{h : \mathcal{X} \rightarrow \{1, 2\}\}$  as the hypothesis space for binary classification, where  
 101 1 represents the ID data, and 2 represents the OOD data. The function  $h$  is called the hypothesis  
 102 function. A scoring function space is a subset of function space, *i.e.*,  $\mathcal{F}_l \subset \{f : \mathcal{X} \rightarrow \mathbb{R}^l\}$ , where  $l$  is  
 103 the output’s dimension of the vector-valued function  $f$ . The function  $f$  is called the scoring function.

104 **Loss and Risks.** Let  $\mathcal{Y}_{\text{all}} = \mathcal{Y} \cup \{K + 1\}$ . Given a loss function  $\ell : \mathcal{Y}_{\text{all}} \times \mathcal{Y}_{\text{all}} \rightarrow \mathbb{R}_{\geq 0}$  satisfying  
 105 that  $\ell(y_1, y_2) = 0$  if and only if  $y_1 = y_2$ , and any  $h \in \mathcal{H}$ , then the *risk* with respect to  $\bar{D}_{XY}$  is

$$R_D(h) := \mathbb{E}_{(\mathbf{x}, y) \sim D_{XY}} \ell(h(\mathbf{x}), y). \quad (1)$$

106 The  $\alpha$ -risk  $R_D^\alpha(h) := (1 - \alpha)R_D^{\text{in}}(h) + \alpha R_D^{\text{out}}(h)$ ,  $\forall \alpha \in [0, 1]$ , where the risks  $R_D^{\text{in}}(h)$ ,  $R_D^{\text{out}}(h)$  are

$$R_D^{\text{in}}(h) := \mathbb{E}_{(\mathbf{x}, y) \sim D_{X_I Y_I}} \ell(h(\mathbf{x}), y), \quad R_D^{\text{out}}(h) := \mathbb{E}_{\mathbf{x} \sim D_{X_O}} \ell(h(\mathbf{x}), K + 1).$$

107 **Learnability.** We aim to select a hypothesis function  $h \in \mathcal{H}$  with approximately minimal risk, based  
 108 on finite data. Generally, we expect the approximation to get better, with the increase in sample size.  
 109 Algorithms achieving this are said to be consistent. Formally, we introduce the following definition:

110 **Definition 1** (Learnability of OOD Detection). *Given a domain space  $\mathcal{D}_{XY}$  and a hypothesis space*  
 111  *$\mathcal{H} \subset \{h : \mathcal{X} \rightarrow \mathcal{Y}_{\text{all}}\}$ , we say OOD detection is **learnable** in  $\mathcal{D}_{XY}$  for  $\mathcal{H}$ , if there exists an algorithm*  
 112  *$\mathbf{A} : \cup_{n=1}^{+\infty} (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathcal{H}$  and a monotonically decreasing sequence  $\epsilon_{\text{cons}}(n)$ , such that  $\epsilon_{\text{cons}}(n) \rightarrow 0$ ,*  
 113 *as  $n \rightarrow +\infty$ , and for any domain  $D_{XY} \in \mathcal{D}_{XY}$ ,*

$$\mathbb{E}_{S \sim D_{X_I Y_I}^n} [R_D(\mathbf{A}(S)) - \inf_{h \in \mathcal{H}} R_D(h)] \leq \epsilon_{\text{cons}}(n), \quad (2)$$

114 *An algorithm  $\mathbf{A}$  for which this holds is said to be consistent with respect to  $\mathcal{D}_{XY}$ .*

115 Definition 1 is a natural extension of agnostic PAC learnability of supervised learning [30]. If for any  
 116  $D_{XY} \in \mathcal{D}_{XY}$ ,  $\pi^{\text{out}} = 0$ , then Definition 2 is the agnostic PAC learnability of supervised learning.  
 117 Although the mathematical expression of Definition 1 is different from the normal definition of  
 118 agnostic PAC learning in [21], one can easily prove that they are equivalent, see Appendix D.3.

119 Since OOD data are unavailable, it is impossible to obtain information about the class-prior probability  
 120  $\pi^{\text{out}}$ . Furthermore, in the real world, it is possible that  $\pi^{\text{out}}$  can be any value in  $[0, 1]$ . Therefore,  
 121 the imbalance issue between ID and OOD distributions, and the priori-unknown issue (*i.e.*,  $\pi^{\text{out}}$  is  
 122 unknown) are the core challenges. To ease these challenges, researchers use AUROC, AUPR and  
 123 FPR95 to estimate the performance of OOD detection [18, 31, 32, 33, 34]. It seems that there is a  
 124 gap between Definition 1 and existing works. To eliminate this gap, we revise Eq. (2) as follows:

$$\mathbb{E}_{S \sim D_{X_I Y_I}^n} [R_D^\alpha(\mathbf{A}(S)) - \inf_{h \in \mathcal{H}} R_D^\alpha(h)] \leq \epsilon_{\text{cons}}(n), \quad \forall \alpha \in [0, 1]. \quad (3)$$

125 If an algorithm  $\mathbf{A}$  satisfies Eq. (3), then the imbalance issue and the prior-unknown issue disappear.  
 126 That is,  $\mathbf{A}$  can simultaneously classify the ID data and detect the OOD data well. Based on the above  
 127 discussion, we define the strong learnability of OOD detection as follows:

128 **Definition 2** (Strong Learnability of OOD Detection). *Given a domain space  $\mathcal{D}_{XY}$  and a hypothesis*  
 129 *space  $\mathcal{H} \subset \{h : \mathcal{X} \rightarrow \mathcal{Y}_{\text{all}}\}$ , we say OOD detection is **strongly learnable** in  $\mathcal{D}_{XY}$  for  $\mathcal{H}$ , if there*  
 130 *exists an algorithm  $\mathbf{A} : \cup_{n=1}^{+\infty} (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathcal{H}$  and a monotonically decreasing sequence  $\epsilon_{\text{cons}}(n)$ ,*  
 131 *such that  $\epsilon_{\text{cons}}(n) \rightarrow 0$ , as  $n \rightarrow +\infty$ , and for any domain  $D_{XY} \in \mathcal{D}_{XY}$ ,*

$$\mathbb{E}_{S \sim D_{X_I Y_I}^n} [R_D^\alpha(\mathbf{A}(S)) - \inf_{h \in \mathcal{H}} R_D^\alpha(h)] \leq \epsilon_{\text{cons}}(n), \quad \forall \alpha \in [0, 1].$$

132 In Theorem 1, we have shown that the strong learnability of OOD detection is equivalent to the  
 133 learnability of OOD detection, if the domain space  $\mathcal{D}_{XY}$  is a *prior-unknown space* (see Definition 3).  
 134 In this paper, we mainly discuss the learnability in the prior-unknown space. Therefore, *when we*  
 135 *mention that OOD detection is learnable, we also mean that OOD detection is strongly learnable.*

136 **Goal of Theory.** Note that the agnostic PAC learnability of supervised learning is distribution-free,  
 137 *i.e.*, the domain space  $\mathcal{D}_{XY}$  consists of all domains. However, due to the absence of OOD data  
 138 during the training process [8, 14, 24], it is obvious that the learnability of OOD detection is not  
 139 distribution-free (*i.e.*, Theorem 4). In fact, we discover that the learnability of OOD detection is  
 140 deeply correlated with the relationship between the domain space  $\mathcal{D}_{XY}$  and the hypothesis space  $\mathcal{H}$ .  
 141 That is, OOD detection is learnable only when the domain space  $\mathcal{D}_{XY}$  and the hypothesis space  $\mathcal{H}$   
 142 satisfy some special conditions, *e.g.*, Condition 1 and Condition 3. We present our goal as follows:

143 **Goal:** *given a hypothesis space  $\mathcal{H}$  and several representative domain spaces  $\mathcal{D}_{XY}$ ,  
 what are the **conditions** to ensure the learnability of OOD detection? Furthermore, if  
 possible, we hope that these conditions are **necessary and sufficient** in some scenarios.*

144 Therefore, compared to the agnostic PAC learnability of supervised learning, our theory doesn't  
 145 focus on the distribution-free case, but focuses on discovering essential conditions to guarantee the  
 146 learnability of OOD detection in several representative and practical domain spaces  $\mathcal{D}_{XY}$ . By these  
 147 essential conditions, we can know *when* OOD detection can be successful in real applications.  
 148 *The guidance for real applications based on our theory and all proofs can be found in Appendices.*

### 149 3 Learning in Priori-unknown Spaces

150 We first investigate a special space, called prior-unknown space. In such space, Definition 1 and  
 151 Definition 2 are equivalent. Furthermore, we also prove that if OOD detection is strongly learnable  
 152 in a space  $\mathcal{D}_{XY}$ , then one can discover a larger domain space, which is prior-unknown, to ensure  
 153 the learnability of OOD detection. These results imply that it is enough to consider our theory in the  
 154 prior-unknown spaces. The prior-unknown space is introduced as follows:

155 **Definition 3.** *Given a domain space  $\mathcal{D}_{XY}$ , we say  $\mathcal{D}_{XY}$  is a priori-unknown space, if for any domain*  
 156  *$D_{XY} \in \mathcal{D}_{XY}$  and any  $\alpha \in [0, 1]$ , we have  $D_{XY}^\alpha := (1 - \alpha)D_{X_1Y_1} + \alpha D_{X_0Y_0} \in \mathcal{D}_{XY}$ .*

157 **Theorem 1.** *Given domain spaces  $\mathcal{D}_{XY}$  and  $\mathcal{D}'_{XY} = \{D_{XY}^\alpha : \forall D_{XY} \in \mathcal{D}_{XY}, \forall \alpha \in [0, 1]\}$ , then*  
 158 *1)  $\mathcal{D}'_{XY}$  is a priori-unknown space and  $\mathcal{D}_{XY} \subset \mathcal{D}'_{XY}$ ;*  
 159 *2) if  $\mathcal{D}_{XY}$  is a priori-unknown space, then Definition 1 and Definition 2 are **equivalent**;*  
 160 *3) OOD detection is strongly learnable in  $\mathcal{D}_{XY}$  **if and only if** OOD detection is learnable in  $\mathcal{D}'_{XY}$ .*

161 The second result of Theorem 1 bridges the learnability and strong learnability, which implies that  
 162 if an algorithm **A** is consistent with respect to a prior-unknown space, then this algorithm **A** can  
 163 address the imbalance issue between ID and OOD distributions, and the priori-unknown issue well.  
 164 Based on Theorem 1, we focus on our theory in the prior-unknown spaces. Furthermore, to demystify  
 165 the learnability of OOD detection, we introduce five representative priori-unknown spaces:

- 166 • Single-distribution space  $\mathcal{D}_{XY}^{D_{XY}}$ . For a domain  $D_{XY}$ ,  $\mathcal{D}_{XY}^{D_{XY}} := \{D_{XY}^\alpha : \forall \alpha \in [0, 1]\}$ .
- 167 • Total space  $\mathcal{D}_{XY}^{\text{all}}$ , which consists of all domains.
- 168 • Separate space  $\mathcal{D}_{XY}^s$ , which consists of all domains that satisfy the separate condition, that is for  
 169 any  $D_{XY} \in \mathcal{D}_{XY}^s$ ,  $\text{supp}D_{X_0} \cap \text{supp}D_{X_1} = \emptyset$ , where  $\text{supp}$  means the support set.
- 170 • Finite-ID-distribution space  $\mathcal{D}_{XY}^F$ , which is a prior-unknown space satisfying that the number of  
 171 distinct ID joint distributions  $D_{X_1Y_1}$  in  $\mathcal{D}_{XY}^F$  is finite, *i.e.*,  $|\{D_{X_1Y_1} : \forall D_{XY} \in \mathcal{D}_{XY}^F\}| < +\infty$ .
- 172 • Density-based space  $\mathcal{D}_{XY}^{\mu, b}$ , which is a prior-unknown space consisting of some domains satisfying  
 173 that: for any  $D_{XY}$ , there exists a density function  $f$  with  $1/b \leq f \leq b$  in  $\text{supp}\mu$  and  $0.5 * D_{X_1} +$   
 174  $0.5 * D_{X_0} = \int f d\mu$ , where  $\mu$  is a measure defined over  $\mathcal{X}$ . Note that if  $\mu$  is discrete, then  $D_X$  is a  
 175 discrete distribution; and if  $\mu$  is the Lebesgue measure, then  $D_X$  is a continuous distribution.

176 The above representative spaces widely exist in real applications. For example, 1) if the images from  
 177 different semantic labels are clearly different (*e.g.*, cats and airplanes), then those images can form  
 178 a distribution belonging to a separate space  $\mathcal{D}_{XY}^s$ ; and 2) when designing an algorithm, we only  
 179 have finite ID datasets, *e.g.*, CIFAR-10, MNIST, SVHN, and ImageNet, to build a model. Then,

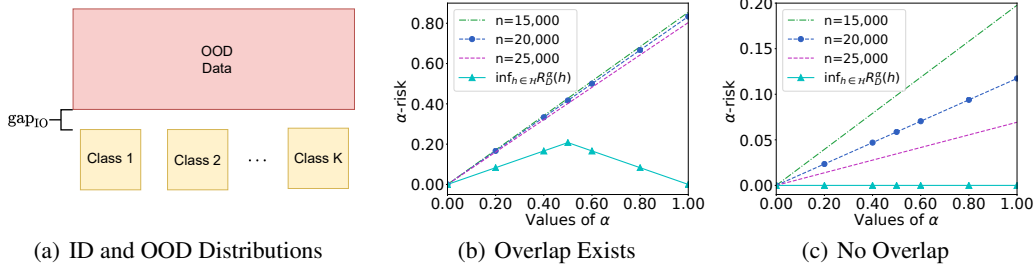


Figure 1: Illustration of  $\inf_{h \in \mathcal{H}} R_D^\alpha(h)$  (solid lines with triangle marks) and the estimated  $\mathbb{E}_{S \sim D_{\text{in}}^n} R_D^\alpha(\mathbf{A}(S))$  (dash lines) with  $\alpha \in [0, 1]$  in different scenarios, where  $D_{\text{in}} = D_{X_I Y_I}$  and the algorithm  $\mathbf{A}$  is the free-energy OOD detection method [23]. Subfigure (a) shows the ID and OOD distributions. In (a),  $\text{gap}_{\text{IO}}$  represents the distance between the support sets of ID and OOD distributions. In (b), since there is an overlap between ID and OOD data, the solid line is a ployline. In (c), since there is no overlap between ID and OOD data, we can check that  $\inf_{h \in \mathcal{H}} R_D^\alpha(h)$  forms a straight line (the solid line). However, since dash lines are always straight lines, two observations can be obtained from (b) and (c): 1) dash lines cannot approximate the solid ployline in (b), which implies the unlearnability of OOD detection; and 2) the solid line in (c) is a straight line and may be approximated by the dash lines in (c). The above observations motivate us to propose Condition 1.

180 finite-ID-distribution space  $\mathcal{D}_{XY}^F$  can handle this real scenario. Note that the single-distribution space  
 181 is a special case of the finite-ID-distribution space. In this paper, we mainly discuss these five spaces.

## 182 4 Impossibility Theorems for OOD Detection

183 In this section, we first give a necessary condition for the learnability of OOD detection. Then, we  
 184 show this necessary condition does not hold in the total space  $\mathcal{D}_{XY}^{\text{all}}$  and the separate space  $\mathcal{D}_{XY}^s$ .

185 **Necessary Condition.** We find a necessary condition for the learnability of OOD detection, *i.e.*, Con-  
 186 dition 1, motivated by the experiments in Figure 1. Details of Figure 1 can be found in Appendix C.3.

187 **Condition 1 (Linear Condition).** For any  $D_{XY} \in \mathcal{D}_{XY}$  and any  $\alpha \in [0, 1]$ ,

$$\inf_{h \in \mathcal{H}} R_D^\alpha(h) = (1 - \alpha) \inf_{h \in \mathcal{H}} R_D^{\text{in}}(h) + \alpha \inf_{h \in \mathcal{H}} R_D^{\text{out}}(h).$$

188 To reveal the importance of Condition 1, Theorem 2 shows that Condition 1 is a *necessary and*  
 189 *sufficient* condition for the learnability of OOD detection if the  $\mathcal{D}_{XY}$  is the single-distribution space.

190 **Theorem 2.** Given a hypothesis space  $\mathcal{H}$  and a domain  $D_{XY}$ , OOD detection is learnable in the  
 single-distribution space  $\mathcal{D}_{XY}^{D_{XY}}$  for  $\mathcal{H}$  if and only if linear condition (*i.e.*, Condition 1) holds.

191 Theorem 2 implies that Condition 1 is important for the learnability of OOD detection. Due to the  
 192 simplicity of single-distribution space, Theorem 2 implies that Condition 1 is the necessary condition  
 193 for the learnability of OOD detection in the prior-unknown space, see Lemma 1 in Appendix F.

194 **Impossibility Theorems.** Here, we first study whether Condition 1 holds in the total space  $\mathcal{D}_{XY}^{\text{all}}$ . If  
 195 Condition 1 does not hold, then OOD detection is not learnable. Theorem 3 shows that Condition 1 is  
 196 not always satisfied, especially, when there is an overlap between the ID and OOD distributions:

197 **Definition 4 (Overlap Between ID and OOD).** We say a domain  $D_{XY}$  has overlap between ID and  
 198 OOD distributions, if there is a  $\sigma$ -finite measure  $\tilde{\mu}$  such that  $D_X$  is absolutely continuous with respect  
 199 to  $\tilde{\mu}$ , and  $\tilde{\mu}(A_{\text{overlap}}) > 0$ , where  $A_{\text{overlap}} = \{\mathbf{x} \in \mathcal{X} : f_1(\mathbf{x}) > 0 \text{ and } f_0(\mathbf{x}) > 0\}$ . Here  $f_1$  and  
 200  $f_0$  are the representers of  $D_{X_I}$  and  $D_{X_O}$  in Radon-Nikodym Theorem [35],

$$D_{X_I} = \int f_1 d\tilde{\mu}, \quad D_{X_O} = \int f_0 d\tilde{\mu}.$$

201 **Theorem 3.** Given a hypothesis space  $\mathcal{H}$  and a prior-unknown space  $\mathcal{D}_{XY}$ , if there is  $D_{XY} \in \mathcal{D}_{XY}$ ,  
 202 which has overlap between ID and OOD, and  $\inf_{h \in \mathcal{H}} R_D^{\text{in}}(h) = 0$  and  $\inf_{h \in \mathcal{H}} R_D^{\text{out}}(h) = 0$ , then  
 203 Condition 1 does not hold. Therefore, OOD detection is not learnable in  $\mathcal{D}_{XY}$  for  $\mathcal{H}$ .

204 Theorem 3 clearly shows that under proper conditions, Condition 1 does not hold, if there exists a  
 205 domain whose ID and OOD distributions have overlap. By Theorem 3, we can obtain that the OOD  
 206 detection is not learnable in the total space  $\mathcal{D}_{XY}^{\text{all}}$  for any non-trivial hypothesis space  $\mathcal{H}$ .

**Theorem 4** (Impossibility Theorem for Total Space). *OOD detection is not learnable in the total space  $\mathcal{D}_{XY}^{\text{all}}$  for  $\mathcal{H}$ , if  $|\phi \circ \mathcal{H}| > 1$ , where  $\phi$  maps ID labels to 1 and maps OOD labels to 2.*

Since the overlaps between ID and OOD distributions may cause that Condition 1 does not hold, we then consider studying the learnability of OOD detection in the separate space  $\mathcal{D}_{XY}^s$ , where there are no overlaps between the ID and OOD distributions. However, Theorem 5 shows that even if we consider the separate space, the OOD detection is still not learnable in some scenarios. Before introducing the impossibility theorem for separate space, *i.e.*, Theorem 5, we need a mild assumption:

**Assumption 1** (Separate Space for OOD). *A hypothesis space  $\mathcal{H}$  is separate for OOD data, if for each data point  $\mathbf{x} \in \mathcal{X}$ , there exists at least one hypothesis function  $h_{\mathbf{x}} \in \mathcal{H}$  such that  $h_{\mathbf{x}}(\mathbf{x}) = K + 1$ .*

Assumption 1 means that every data point  $\mathbf{x}$  has the possibility to be detected as OOD data. Assumption 1 is mild and can be satisfied by many hypothesis spaces, *e.g.*, the FCNN-based hypothesis space (Proposition 1 in Appendix K), score-based hypothesis space (Proposition 2 in Appendix K) and universal kernel space. Next, we use *Vapnik–Chervonenkis* (VC) dimension [22] to measure the size of hypothesis space, and study the learnability of OOD detection in  $\mathcal{D}_{XY}^s$  based on the VC dimension.

**Theorem 5** (Impossibility Theorem for Separate Space). *If Assumption 1 holds,  $\text{VCdim}(\phi \circ \mathcal{H}) < +\infty$  and  $\sup_{h \in \mathcal{H}} |\{\mathbf{x} \in \mathcal{X} : h(\mathbf{x}) \in \mathcal{Y}\}| = +\infty$ , then OOD detection is not learnable in separate space  $\mathcal{D}_{XY}^s$  for  $\mathcal{H}$ , where  $\phi$  maps ID labels to 1 and maps OOD labels to 2.*

The finite VC dimension normally implies the learnability of supervised learning. However, in our results, the finite VC dimension cannot guarantee the learnability of OOD detection in the separate space, which reveals the difficulty of the OOD detection. Although the above impossibility theorems are frustrating, there is still room to discuss the conditions in Theorem 5, and to find out the proper conditions for ensuring the learnability of OOD detection in the separate space (see Sections 5 and 6).

## 5 When OOD Detection Can Be Successful

Here, we discuss when the OOD detection can be successful in the separate space  $\mathcal{D}_{XY}^s$ , finite-ID-distribution space  $\mathcal{D}_{XY}^F$  and density-based space  $\mathcal{D}_{XY}^{\mu, b}$ . We first study the separate space  $\mathcal{D}_{XY}^s$ .

**OOD Detection in the Separate Space.** Theorem 5 has indicated that  $\text{VCdim}(\phi \circ \mathcal{H}) = +\infty$  or  $\sup_{h \in \mathcal{H}} |\{\mathbf{x} \in \mathcal{X} : h(\mathbf{x}) \in \mathcal{Y}\}| < +\infty$  is necessary to ensure the learnability of OOD detection in  $\mathcal{D}_{XY}^s$  if Assumption 1 holds. However, generally, hypothesis spaces generated by feed-forward neural networks with proper activation functions have finite VC dimension [36, 37]. Therefore, we study the learnability of OOD detection in the case that  $|\mathcal{X}| < +\infty$ , which implies that  $\sup_{h \in \mathcal{H}} |\{\mathbf{x} \in \mathcal{X} : h(\mathbf{x}) \in \mathcal{Y}\}| < +\infty$ . Additionally, Theorem 10 also implies that  $|\mathcal{X}| < +\infty$  is the necessary and sufficient condition for the learnability of OOD detection in separate space, when the hypothesis space is generated by FCNN. Hence,  $|\mathcal{X}| < +\infty$  may be necessary in the space  $\mathcal{D}_{XY}^s$ .

For simplicity, we first discuss the case that  $K = 1$ , *i.e.*, the one-class novelty detection. We show the necessary and sufficient condition for the learnability of OOD detection in  $\mathcal{D}_{XY}^s$ , when  $|\mathcal{X}| < +\infty$ .

**Theorem 6.** *Let  $K = 1$  and  $|\mathcal{X}| < +\infty$ . Suppose that Assumption 1 holds and the constant function  $h^{\text{in}} := 1 \in \mathcal{H}$ . Then OOD detection is learnable in  $\mathcal{D}_{XY}^s$  for  $\mathcal{H}$  if and only if  $\mathcal{H}_{\text{all}} - \{h^{\text{out}}\} \subset \mathcal{H}$ , where  $\mathcal{H}_{\text{all}}$  is the hypothesis space consisting of all hypothesis functions, and  $h^{\text{out}}$  is a constant function that  $h^{\text{out}} := 2$ , here 1 represents ID data and 2 represents OOD data.*

The condition  $h^{\text{in}} \in \mathcal{H}$  presented in Theorem 6 is mild. Many practical hypothesis spaces satisfy this condition, *e.g.*, the FCNN-based hypothesis space (Proposition 1 in Appendix K), score-based hypothesis space (Proposition 2 in Appendix K) and universal kernel-based hypothesis space. Theorem 6 implies that if  $K = 1$  and OOD detection is learnable in  $\mathcal{D}_{XY}^s$  for  $\mathcal{H}$ , then the hypothesis space  $\mathcal{H}$  should contain almost all hypothesis functions, implying that if the OOD detection can be learnable in the distribution-agnostic case, then a large-capacity model is necessary.

Next, we extend Theorem 6 to a general case, *i.e.*,  $K > 1$ . When  $K > 1$ , we will first use a binary classifier  $h^b$  to classify the ID and OOD data. Then, for the ID data identified by  $h^b$ , an ID hypothesis function  $h^{\text{in}}$  will be used to classify them into corresponding ID classes. We state this strategy as follows: given a hypothesis space  $\mathcal{H}^{\text{in}}$  for ID distribution and a binary classification hypothesis space  $\mathcal{H}^b$  introduced in Section 2, we use  $\mathcal{H}^{\text{in}}$  and  $\mathcal{H}^b$  to construct an OOD detection’s hypothesis space  $\mathcal{H}$ , which consists of all hypothesis functions  $h$  satisfying the following condition: there exist  $h^{\text{in}} \in \mathcal{H}^{\text{in}}$

254 and  $h^b \in \mathcal{H}^b$  such that for any  $\mathbf{x} \in \mathcal{X}$ ,

$$h(\mathbf{x}) = i, \quad \text{if } h^{\text{in}}(\mathbf{x}) = i \text{ and } h^b(\mathbf{x}) = 1; \text{ otherwise, } h(\mathbf{x}) = K + 1. \quad (4)$$

255 We use  $\mathcal{H}^{\text{in}} \bullet \mathcal{H}^b$  to represent a hypothesis space consisting of all  $h$  defined in Eq. (4). In addition,  
256 we also need an additional condition for the loss function  $\ell$ . This condition is shown as follows:

257 **Condition 2.**  $\ell(y_2, y_1) \leq \ell(K + 1, y_1)$ , for any in-distribution labels  $y_1$  and  $y_2 \in \mathcal{Y}$ .

258 **Theorem 7.** Let  $|\mathcal{X}| < +\infty$  and  $\mathcal{H} = \mathcal{H}^{\text{in}} \bullet \mathcal{H}^b$ . If  $\mathcal{H}_{\text{all}} - \{h^{\text{out}}\} \subset \mathcal{H}^b$  and Condition 2 holds,  
259 then OOD detection is learnable in  $\mathcal{D}_{XY}^s$  for  $\mathcal{H}$ , where  $\mathcal{H}_{\text{all}}$  and  $h^{\text{out}}$  are defined in Theorem 6.

260 **OOD Detection in the Finite-ID-Distribution Space.** Since researchers can only collect finite ID  
261 datasets as the training data in the process of algorithm design, it is worthy to study the learnability of  
262 OOD detection in the finite-ID-distribution space  $\mathcal{D}_{XY}^F$ . We first show two necessary concepts below.

263 **Definition 5 (ID Consistency).** Given a domain space  $\mathcal{D}_{XY}$ , we say any two domains  $D_{XY} \in \mathcal{D}_{XY}$   
264 and  $D'_{XY} \in \mathcal{D}_{XY}$  are ID consistency, if  $D_{X_1Y_1} = D'_{X_1Y_1}$ . We use the notation  $\sim$  to represent the ID  
265 consistency, i.e.,  $D_{XY} \sim D'_{XY}$  if and only if  $D_{XY}$  and  $D'_{XY}$  are ID consistency.

266 It is easy to check that the ID consistency  $\sim$  is an equivalence relation. Therefore, we define the set  
267  $[D_{XY}] := \{D'_{XY} \in \mathcal{D}_{XY} : D_{XY} \sim D'_{XY}\}$  as the equivalence class with respect to space  $\mathcal{D}_{XY}$ .

268 **Condition 3 (Compatibility).** For any equivalence class  $[D'_{XY}]$  with respect to  $\mathcal{D}_{XY}$  and any  $\epsilon > 0$ ,  
269 there exists a hypothesis function  $h_\epsilon \in \mathcal{H}$  such that for any domain  $D_{XY} \in [D'_{XY}]$ ,

$$h_\epsilon \in \{h' \in \mathcal{H} : R_D^{\text{out}}(h') \leq \inf_{h \in \mathcal{H}} R_D^{\text{out}}(h) + \epsilon\} \cap \{h' \in \mathcal{H} : R_D^{\text{in}}(h') \leq \inf_{h \in \mathcal{H}} R_D^{\text{in}}(h) + \epsilon\}.$$

270 In Appendix F, Lemma 2 has implied that Condition 3 is a general version of Condition 1. Next,  
271 Theorem 8 indicates that Condition 3 is the necessary and sufficient condition in the space  $\mathcal{D}_{XY}^F$ .

272 **Theorem 8.** Suppose that  $\mathcal{X}$  is a bounded set. OOD detection is learnable in the finite-ID-  
distribution space  $\mathcal{D}_{XY}^F$  for  $\mathcal{H}$  if and only if the compatibility condition (i.e., Condition 3) holds.  
Furthermore, the learning rate  $\epsilon_{\text{cons}}(n)$  can attain  $O(1/\sqrt{n^{1-\theta}})$ , for any  $\theta \in (0, 1)$ .

273 Theorem 8 shows that, in the process of algorithm design, OOD detection cannot be successful  
274 without the compatibility condition. Theorem 8 also implies that Condition 3 is essential for the  
275 learnability of OOD detection. This motivates us to study whether OOD detection can be successful  
276 in more general spaces (e.g., the density-based space), when the compatibility condition holds.

277 **OOD Detection in the Density-based Space.** To ensure that Condition 3 holds, we consider a basic  
278 assumption in learning theory—*Realizability Assumption*, i.e., for any  $D_{XY} \in \mathcal{D}_{XY}$ , there exists  
279  $h^* \in \mathcal{H}$  such that  $R_D(h^*) = 0$ . We discover that in the density-based space  $\mathcal{D}_{XY}^{\mu, b}$ , Realizability  
280 Assumption can conclude the compatibility condition (i.e., Condition 3). Based on this observation,  
281 we can prove the following theorem:

282 **Theorem 9.** Given a density-based space  $\mathcal{D}_{XY}^{\mu, b}$ , if  $\mu(\mathcal{X}) < +\infty$ , the Realizability Assumption  
holds, then when  $\mathcal{H}$  has finite Natarajan dimension [21], OOD detection is learnable in  $\mathcal{D}_{XY}^{\mu, b}$  for  
 $\mathcal{H}$ . Furthermore, the learning rate  $\epsilon_{\text{cons}}(n)$  can attain  $O(1/\sqrt{n^{1-\theta}})$ , for any  $\theta \in (0, 1)$ .

283 To further investigate the importance and necessary of Realizability Assumption, Theorem 11 has  
284 indicated that in some practical scenarios, Realizability Assumption is the necessary and sufficient  
285 condition for the learnability of OOD detection in the density-based space. Therefore, Realizability  
286 Assumption may be indispensable for the learnability of OOD detection in some practical scenarios.

## 287 6 Connecting Theory to Practice

288 In Section 5, we have shown the successful scenarios where OOD detection problem can be addressed  
289 in theory. In this section, we will discuss how the proposed theory is applied to two representative  
290 hypothesis spaces—neural-network-based hypothesis spaces and score-based hypothesis spaces.

291 **Fully-connected Neural Networks.** Given a sequence  $\mathbf{q} = (l_1, l_2, \dots, l_g)$ , where  $l_i$  and  $g$  are positive  
292 integers and  $g > 2$ , we use  $g$  to represent the *depth* of neural network and use  $l_i$  to represent the *width*

293 of the  $i$ -th layer. After the activation function  $\sigma$  is selected<sup>1</sup>, we can obtain the architecture of FCNN  
 294 according to the sequence  $\mathbf{q}$ . Let  $\mathbf{f}_{\mathbf{w},\mathbf{b}}$  be the function generated by FCNN with weights  $\mathbf{w}$  and bias  
 295  $\mathbf{b}$ . An FCNN-based scoring function space is defined as:  $\mathcal{F}_{\mathbf{q}}^{\sigma} := \{\mathbf{f}_{\mathbf{w},\mathbf{b}} : \forall \text{ weights } \mathbf{w}, \forall \text{ bias } \mathbf{b}\}$ . In  
 296 addition, for simplicity, given any two sequences  $\mathbf{q} = (l_1, \dots, l_g)$  and  $\mathbf{q}' = (l'_1, \dots, l'_{g'})$ , we use the  
 297 notation  $\mathbf{q} \lesssim \mathbf{q}'$  to represent the following equations and inequalities:

$$1) g \leq g', l_1 = l'_1, l_g = l'_{g'}; \quad 2) l_i \leq l'_i, \forall i = 1, \dots, g-1; \quad \text{and} \quad 3) l_{g-1} \leq l'_i, \forall i = g, \dots, g'-1.$$

298 In Appendix L, Lemma 10 shows  $\mathbf{q} \lesssim \mathbf{q}' \Rightarrow \mathcal{F}_{\mathbf{q}}^{\sigma} \subset \mathcal{F}_{\mathbf{q}'}^{\sigma}$ . We use  $\lesssim$  to compare the sizes of FCNNs.

299 **FCNN-based Hypothesis Space.** Let  $l_g = K + 1$ . The FCNN-based scoring function space  $\mathcal{F}_{\mathbf{q}}^{\sigma}$  can  
 300 induce an FCNN-based hypothesis space. For any  $\mathbf{f}_{\mathbf{w},\mathbf{b}} \in \mathcal{F}_{\mathbf{q}}^{\sigma}$ , the induced hypothesis function is:

$$h_{\mathbf{w},\mathbf{b}} := \arg \max_{k \in \{1, \dots, K+1\}} f_{\mathbf{w},\mathbf{b}}^k, \text{ where } f_{\mathbf{w},\mathbf{b}}^k \text{ is the } k\text{-th coordinate of } \mathbf{f}_{\mathbf{w},\mathbf{b}}.$$

301 Then, the FCNN-based hypothesis space is defined as  $\mathcal{H}_{\mathbf{q}}^{\sigma} := \{h_{\mathbf{w},\mathbf{b}} : \forall \text{ weights } \mathbf{w}, \forall \text{ bias } \mathbf{b}\}$ .

302 **Score-based Hypothesis Space.** Many OOD detection algorithms detect OOD data by using a  
 303 score-based strategy. That is, given a threshold  $\lambda$ , a scoring function space  $\mathcal{F}_l \subset \{\mathbf{f} : \mathcal{X} \rightarrow \mathbb{R}^l\}$   
 304 and a scoring function  $E : \mathcal{F}_l \rightarrow \mathbb{R}$ , then  $\mathbf{x}$  is regarded as ID data if and only if  $E(\mathbf{f}(\mathbf{x})) \geq \lambda$ . We  
 305 introduce several representative scoring functions  $E$  as follows: for any  $\mathbf{f} = [f^1, \dots, f^l]^{\top} \in \mathcal{F}_l$ ,

306 • softmax-based function [7] and temperature-scaled function [8]:  $\lambda \in (\frac{1}{l}, 1)$  and  $T > 0$ ,

$$E(\mathbf{f}) = \max_{k \in \{1, \dots, l\}} \frac{\exp(f^k)}{\sum_{c=1}^l \exp(f^c)}, \quad E(\mathbf{f}) = \max_{k \in \{1, \dots, l\}} \frac{\exp(f^k/T)}{\sum_{c=1}^l \exp(f^c/T)}; \quad (5)$$

307 • energy-based function [23]:  $\lambda \in (0, +\infty)$  and  $T > 0$ ,

$$E(\mathbf{f}) = T \log \sum_{c=1}^l \exp(f^c/T). \quad (6)$$

308 Using  $E$ ,  $\lambda$  and  $\mathbf{f} \in \mathcal{F}_{\mathbf{q}}^{\sigma}$ , we have a classifier:  $h_{\mathbf{f},E}^{\lambda}(\mathbf{x}) = 1$ , if  $E(\mathbf{f}(\mathbf{x})) \geq \lambda$ ; otherwise,  $h_{\mathbf{f},E}^{\lambda}(\mathbf{x}) = 2$ ,  
 309 where 1 represents the ID data and 2 represents the OOD data. Hence, a binary classification  
 310 hypothesis space  $\mathcal{H}^b$ , which consists of all  $h_{\mathbf{f},E}^{\lambda}$ , is generated. We define  $\mathcal{H}_{\mathbf{q},E}^{\sigma,\lambda} := \{h_{\mathbf{f},E}^{\lambda} : \forall \mathbf{f} \in \mathcal{F}_{\mathbf{q}}^{\sigma}\}$ .

311 **Learnability of OOD Detection in Different Hypothesis Spaces.** Next, we present applications of  
 312 our theory regarding the above two practical and important hypothesis spaces  $\mathcal{H}_{\mathbf{q}}^{\sigma}$  and  $\mathcal{H}_{\mathbf{q},E}^{\sigma,\lambda}$ .

313 **Theorem 10.** *Suppose that Condition 2 holds and the hypothesis space  $\mathcal{H}$  is FCNN-based or score-*  
 314 *based, i.e.,  $\mathcal{H} = \mathcal{H}_{\mathbf{q}}^{\sigma}$  or  $\mathcal{H} = \mathcal{H}^{\text{in}} \bullet \mathcal{H}^b$ , where  $\mathcal{H}^{\text{in}}$  is an ID hypothesis space,  $\mathcal{H}^b = \mathcal{H}_{\mathbf{q},E}^{\sigma,\lambda}$  and*  
 315  *$\mathcal{H} = \mathcal{H}^{\text{in}} \bullet \mathcal{H}^b$  is introduced below Eq. (4), here  $E$  is introduced in Eqs. (5) or (6). Then*

316 *There is a sequence  $\mathbf{q} = (l_1, \dots, l_g)$  such that OOD detection is learnable in the separate space  $\mathcal{D}_{XY}^s$  for  $\mathcal{H}$  if and only if  $|\mathcal{X}| < +\infty$ .*

317 *Furthermore, if  $|\mathcal{X}| < +\infty$ , then there exists a sequence  $\mathbf{q} = (l_1, \dots, l_g)$  such that for any sequence*  
 318  *$\mathbf{q}'$  satisfying that  $\mathbf{q} \lesssim \mathbf{q}'$ , OOD detection is learnable in  $\mathcal{D}_{XY}^s$  for  $\mathcal{H}$ .*

319 Theorem 10 states that 1) when the hypothesis space is FCNN-based or score-based, the finite feature  
 320 space is the necessary and sufficient condition for the learnability of OOD detection in the separate  
 321 space; and 2) a larger architecture of FCNN has a greater probability to achieve the learnability of  
 322 OOD detection in the separate space. Note that when we select Eqs. (5) or (6) as the scoring function  
 323  $E$ , Theorem 10 also shows that the selected scoring functions  $E$  can guarantee the learnability of  
 324 OOD detection, which is a theoretical support for the representative works [8, 23, 7]. Furthermore,  
 325 Theorem 11 also offers theoretical supports for these works in the density-based space, when  $K = 1$ .

326 **Theorem 11.** *Suppose that each domain  $D_{XY}$  in  $\mathcal{D}_{XY}^{\mu,b}$  is attainable, i.e.,  $\arg \min_{h \in \mathcal{H}} R_D(h) \neq \emptyset$*   
 327 *(the finite discrete domains satisfy this). Let  $K = 1$  and the hypothesis space  $\mathcal{H}$  be score-based*  
 328 *( $\mathcal{H} = \mathcal{H}_{\mathbf{q},E}^{\sigma,\lambda}$ , where  $E$  is in Eqs. (5) or (6)) or FCNN-based ( $\mathcal{H} = \mathcal{H}_{\mathbf{q}}^{\sigma}$ ). If  $\mu(\mathcal{X}) < +\infty$ , then the*  
 329 *following four conditions are equivalent:*

330 *Learnability in  $\mathcal{D}_{XY}^{\mu,b}$  for  $\mathcal{H}$   $\iff$  Condition 1  $\iff$  Realizability Assumption  $\iff$  Condition 3*

<sup>1</sup>We consider the *rectified linear unit* (ReLU) function as the default activation function  $\sigma$ , which is defined by  $\sigma(x) = \max\{x, 0\}$ ,  $\forall x \in \mathbb{R}$ . We will not repeatedly mention the definition of  $\sigma$  in the rest of our paper.

331 Theorem 11 still holds if the function space  $\mathcal{F}_q^\sigma$  is generated by Convolutional Neural Network.

332 **Overlap and Benefits of Multi-class Case.** We investigate when the hypothesis space is FCNN-based  
333 or score-based, what will happen if there exists an overlap between the ID and OOD distributions?

334 **Theorem 12.** Let  $K = 1$  and the hypothesis space  $\mathcal{H}$  be score-based ( $\mathcal{H} = \mathcal{H}_{q,E}^{\sigma,\lambda}$ , where  $E$  is in  
335 Eqs. (5) or (6)) or FCNN-based ( $\mathcal{H} = \mathcal{H}_q^\sigma$ ). Given a prior-unknown space  $\mathcal{D}_{XY}$ , if there exists a  
336 domain  $D_{XY} \in \mathcal{D}_{XY}$ , which has an overlap between ID and OOD distributions (see Definition 4),  
337 then OOD detection is not learnable in the domain space  $\mathcal{D}_{XY}$  for  $\mathcal{H}$ .

338 When  $K = 1$  and the hypothesis space is FCNN-based or score-based, Theorem 12 shows that overlap  
339 between ID and OOD distributions is the sufficient condition for the unlearnability of OOD detection.  
340 Theorem 12 takes roots in the conditions  $\inf_{h \in \mathcal{H}} R_D^{\text{in}}(h) = 0$  and  $\inf_{h \in \mathcal{H}} R_D^{\text{out}}(h) = 0$ . However,  
341 when  $K > 1$ , we can ensure  $\inf_{h \in \mathcal{H}} R_D^{\text{in}}(h) > 0$  if ID distribution  $D_{X_1 Y_1}$  has overlap between ID  
342 classes. By this observation, we conjecture that when  $K > 1$ , OOD detection is learnable in some  
343 special cases where overlap exists, even if the hypothesis space is FCNN-based or score-based.

## 344 7 Related Work

345 We briefly review the related theoretical works below. See Appendix A for detailed related works.

346 **OOD Detection Theory.** [38] understands the OOD detection via goodness-of-fit tests and typical  
347 set hypothesis, and argues that minimal density estimation errors can lead to OOD detection failures,  
348 when there exists an overlap between ID and OOD distributions. Beyond [38], [39] paves a new  
349 avenue to designing provable OOD detection algorithms. Compared to [39, 38], our theory focuses  
350 on the agnostic PAC learnable theory of OOD detection and identifies several necessary and sufficient  
351 conditions for the learnability of OOD detection, opening a door to study OOD detection in theory.

352 **Open-set Learning Theory.** [40] and [29, 41] propose the agnostic PAC learning bounds for open-set  
353 detection and open-set domain adaptation, respectively. Unfortunately, [29, 40, 41] all require that  
354 the test data are indispensable during the training process. To investigate open-set learning (OSL)  
355 *without accessing the test data* during training, [24] proposes and investigates the *almost* agnostic  
356 PAC learnability for OSL. However, the assumptions used in [24] are very strong and unpractical.

357 **Learning Theory for Classification with Reject Option.** Many works [42, 43] also investigate the  
358 *classification with reject option* (CwRO) problem, which is similar to OOD detection in some cases.  
359 [44, 45, 46, 47, 48] study the learning theory and propose the agnostic PAC learning bounds for  
360 CwRO. However, compared to our work regarding OOD detection, existing CwRO theories mainly  
361 focus on how the ID risk  $R_D^{\text{in}}$  (*i.e.*, the risk that ID data is wrongly classified) is influenced by special  
362 rejection rules. Our theory not only focuses on the ID risk, but also pays attention to the OOD risk.

363 **PQ Learning Theory.** Under some conditions, PQ learning theory [49, 50] can be regarded as the  
364 PAC theory for OOD detection in the semi-supervised or transductive learning cases, *i.e.*, test data  
365 are required during training. Besides, [49, 50] aim to give the PAC estimation under Realizability  
366 Assumption [21]. Our theory does not only study the PAC estimation in the realization cases, but also  
367 studies the agnostic cases, which are more difficult than PAC theory under Realizability Assumption.

## 368 8 Conclusions

369 Detecting OOD data has shown its significance in improving the reliability of machine learning.  
370 However, very few works discuss OOD detection in theory, which hinders real-world applications  
371 of OOD detection algorithms. In this paper, we are the *first* to provide the agnostic PAC theory for  
372 OOD detection. Our results imply that we cannot expect a universally consistent algorithm to handle  
373 all scenarios in OOD detection. Yet, it is still possible to make OOD detection learnable in certain  
374 scenarios. For example, when the ID and OOD images have clearly different semantic meanings,  
375 Theorems 10 and 11 show that the image-based OOD detection is learnable for some practical  
376 hypothesis spaces. In addition, when we design OOD detection algorithms, we normally only have  
377 finite ID datasets. In this real scenario, Theorem 8 provides a necessary and sufficient condition for  
378 the success of OOD detection. Our theory reveals many necessary and sufficient condition for the  
379 learnability of OOD detection, hence *opening a door* to studying the learnability of OOD detection.

## 380 References

- 381 [1] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
382 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,  
383 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image  
384 recognition at scale. In *ICLR*, 2021.
- 385 [2] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected  
386 convolutional networks. In *CVPR*, 2017.
- 387 [3] Yen-Chang Hsu, Yilin Shen, Hongxia Jin, and Zsolt Kira. Generalized ODIN: detecting  
388 out-of-distribution image without learning from out-of-distribution data. In *CVPR*, 2020.
- 389 [4] Jingkang Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu. Generalized out-of-distribution  
390 detection: A survey. *CoRR*, abs/2110.11334, 2021.
- 391 [5] Abhijit Bendale and Terrance E Boult. Towards open set deep networks. In *The IEEE / CVF*  
392 *Computer Vision and Pattern Recognition Conference (CVPR)*, 2016.
- 393 [6] Jiefeng Chen, Yixuan Li, Xi Wu, Yingyu Liang, and Somesh Jha. Atom: Robustifying out-of-  
394 distribution detection using outlier mining. *ECML*, 2021.
- 395 [7] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution  
396 examples in neural networks. In *ICLR*, 2017.
- 397 [8] Shiyu Liang, Yixuan Li, and R. Srikant. Enhancing the reliability of out-of-distribution image  
398 detection in neural networks. In *ICLR*, 2018.
- 399 [9] Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for  
400 detecting out-of-distribution samples and adversarial attacks. In *NeurIPS*, 2018.
- 401 [10] Bo Zong, Qi Song, Martin Renqiang Min, Wei Cheng, Cristian Lumezanu, Dae-ki Cho, and  
402 Haifeng Chen. Deep autoencoding gaussian mixture model for unsupervised anomaly detection.  
403 In *ICLR*, 2018.
- 404 [11] Stanislav Pidhorskyi, Ranya Almohsen, and Gianfranco Doretto. Generative probabilistic  
405 novelty detection with adversarial autoencoders. In *NeurIPS*, 2018.
- 406 [12] Eric T. Nalisnick, Akihiro Matsukawa, Yee Whye Teh, Dilan Görür, and Balaji Lakshmi-  
407 narayanan. Do deep generative models know what they don't know? In *ICLR*, 2019.
- 408 [13] Dan Hendrycks, Mantas Mazeika, and Thomas G. Dietterich. Deep anomaly detection with  
409 outlier exposure. In *ICLR*, 2019.
- 410 [14] Jie Ren, Peter J. Liu, Emily Fertig, Jasper Snoek, Ryan Poplin, Mark A. DePristo, Joshua V.  
411 Dillon, and Balaji Lakshminarayanan. Likelihood ratios for out-of-distribution detection. In  
412 *NeurIPS*, 2019.
- 413 [15] Ziqian Lin, Sreya Dutta Roy, and Yixuan Li. Mood: Multi-level out-of-distribution detection.  
414 In *CVPR*, 2021.
- 415 [16] Mohammadreza Salehi, Hossein Mirzaei, Dan Hendrycks, Yixuan Li, Mohammad Hossein  
416 Rohban, and Mohammad Sabokrou. A unified survey on anomaly, novelty, open-set, and  
417 out-of-distribution detection: Solutions and future challenges. *arXiv preprint arXiv:2110.14051*,  
418 2021.
- 419 [17] Yiyu Sun, Chuan Guo, and Yixuan Li. React: Out-of-distribution detection with rectified  
420 activations. In *NeurIPS*, 2021.
- 421 [18] Rui Huang, Andrew Geng, and Yixuan Li. On the Importance of Gradients for Detecting  
422 Distributional Shifts in the Wild. In *NeurIPS*, 2021.
- 423 [19] Stanislav Fort, Jie Ren, and Balaji Lakshminarayanan. Exploring the Limits of Out-of-  
424 Distribution Detection. In *NeurIPS*, 2021.

- 425 [20] Yifei Ming, Hang Yin, and Yixuan Li. On the impact of spurious correlation for out-of-  
426 distribution detection. *AAAI*, 2022.
- 427 [21] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to*  
428 *algorithms*. Cambridge university press, 2014.
- 429 [22] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*.  
430 MIT press, 2018.
- 431 [23] Weitang Liu, Xiaoyun Wang, John D. Owens, and Yixuan Li. Energy-based out-of-distribution  
432 detection. In *NeurIPS*, 2020.
- 433 [24] Zhen Fang, Jie Lu, Anjin Liu, Feng Liu, and Guangquan Zhang. Learning bounds for open-set  
434 learning. In *ICML*, 2021.
- 435 [25] Guangyao Chen, Peixi Peng, Xiangqian Wang, and Yonghong Tian. Adversarial reciprocal  
436 points learning for open set recognition. *IEEE Transactions on Pattern Analysis and Machine*  
437 *Intelligence*, 2021.
- 438 [26] Lukas Ruff, Nico Görnitz, Lucas Deecke, Shoaib Ahmed Siddiqui, Robert A. Vandermeulen,  
439 Alexander Binder, Emmanuel Müller, and Marius Kloft. Deep one-class classification. In *ICML*,  
440 2018.
- 441 [27] Sachin Goyal, Aditi Raghunathan, Moksh Jain, Harsha Vardhan Simhadri, and Prateek Jain.  
442 DROCC: deep robust one-class classification. In *ICML*, 2020.
- 443 [28] Lucas Deecke, Robert A. Vandermeulen, Lukas Ruff, Stephan Mandt, and Marius Kloft. Image  
444 anomaly detection with generative adversarial networks. In *ECML*, 2018.
- 445 [29] Z. Fang, Jie Lu, F. Liu, Junyu Xuan, and G. Zhang. Open set domain adaptation: Theoretical  
446 bound and algorithm. *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- 447 [30] Shai Shalev-Shwartz, Ohad Shamir, Nathan Srebro, and Karthik Sridharan. Learnability,  
448 stability and uniform convergence. *J. Mach. Learn. Res.*, 11:2635–2670, 2010.
- 449 [31] Guangyao Chen, Limeng Qiao, Yemin Shi, Peixi Peng, Jia Li, Tiejun Huang, Shiliang Pu, and  
450 Yonghong Tian. Learning open set network with discriminative reciprocal points. *ICCV*, 2020.
- 451 [32] Jiefeng Chen, Yixuan Li, Xi Wu, Yingyu Liang, and Somesh Jha. Informative outlier matters:  
452 Robustifying out-of-distribution detection using outlier mining. *ICML Workshop*, 2020.
- 453 [33] Jiefeng Chen, Yixuan Li, Xi Wu, Yingyu Liang, and Somesh Jha. Robust out-of-distribution  
454 detection for neural networks. *arXiv preprint arXiv:2003.09711*, 2020.
- 455 [34] Wentao Bao, Qi Yu, and Yu Kong. Evidential deep learning for open set action recognition.  
456 *ICCV*, 2021.
- 457 [35] Donald L Cohn. *Measure theory*. Springer, 2013.
- 458 [36] Peter L. Bartlett, Nick Harvey, Christopher Liaw, and Abbas Mehrabian. Nearly-tight vc-  
459 dimension and pseudodimension bounds for piecewise linear neural networks. *Journal of*  
460 *Machine Learning Research*, 20(63):1–17, 2019.
- 461 [37] Marek Karpinski and Angus Macintyre. Polynomial bounds for VC dimension of sigmoidal  
462 and general pfaffian neural networks. *J. Comput. Syst. Sci.*, 54(1):169–176, 1997.
- 463 [38] Lily H. Zhang, Mark Goldstein, and Rajesh Ranganath. Understanding failures in out-of-  
464 distribution detection with deep generative models. In *ICML*, 2021.
- 465 [39] Peyman Morteza and Yixuan Li. Provable guarantees for understanding out-of-distribution  
466 detection. *AAAI*, 2022.
- 467 [40] Si Liu, Risheek Garrepalli, Thomas G. Dietterich, Alan Fern, and Dan Hendrycks. Open  
468 category detection with PAC guarantees. In *ICML*, 2018.

- 469 [41] Yadan Luo, Zijian Wang, Zi Huang, and Mahsa Baktashmotlagh. Progressive graph learning  
470 for open-set domain adaptation. In *ICML*, 2020.
- 471 [42] C. K. Chow. On optimum recognition error and reject tradeoff. *IEEE Transactions on Informa-*  
472 *tion Theory*, 1970.
- 473 [43] Vojtech Franc, Daniel Průša, and V. Voracek. Optimal strategies for reject option classifiers.  
474 *CoRR*, abs/2101.12523, 2021.
- 475 [44] Corinna Cortes, Giulia DeSalvo, and Mehryar Mohri. Learning with rejection. In *ALT*, 2016.
- 476 [45] Corinna Cortes, Giulia DeSalvo, and Mehryar Mohri. Boosting with abstention. In *NeurIPS*,  
477 2016.
- 478 [46] Chenri Ni, Nontawat Charoenphakdee, Junya Honda, and Masashi Sugiyama. On the calibration  
479 of multiclass classification with rejection. In *NeurIPS*, 2019.
- 480 [47] Nontawat Charoenphakdee, Zhenghang Cui, Yivan Zhang, and Masashi Sugiyama. Classifica-  
481 tion with rejection based on cost-sensitive classification. In *ICML*, 2021.
- 482 [48] Peter L. Bartlett and Marten H. Wegkamp. Classification with a reject option using a hinge loss.  
483 *Journal of Machine Learning Research*, 2008.
- 484 [49] Shafi Goldwasser, Adam Tauman Kalai, Yael Kalai, and Omar Montasser. Beyond perturbations:  
485 Learning guarantees with arbitrary adversarial test examples. In *NeurIPS*, 2020.
- 486 [50] Adam Tauman Kalai and Varun Kanade. Efficient learning with arbitrary covariate shift. In  
487 *ALT*, Proceedings of Machine Learning Research, 2021.
- 488 [51] Akshay Raj Dhamija, Manuel Günther, and Terrance E. Boult. Reducing network agnostophobia.  
489 In *NeurIPS*, pages 9175–9186, 2018.
- 490 [52] Haoran Wang, Weitang Liu, Alex Bocchieri, and Yixuan Li. Can multi-label classification  
491 networks know what they don’t know? In *NeurIPS*, 2021.
- 492 [53] Diederik P. Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolu-  
493 tions. In *NeurIPS*, 2018.
- 494 [54] Zhisheng Xiao, Qing Yan, and Yali Amit. Likelihood regret: An out-of-distribution detection  
495 score for variational auto-encoder. In *NeurIPS*, 2020.
- 496 [55] Jie Ren, Stanislav Fort, Jeremiah Liu, Abhijit Guha Roy, Shreyas Padhy, and Balaji Lakshmi-  
497 narayanan. A simple fix to mahalalanobis distance for improving near-ood detection. *CoRR*,  
498 abs/2106.09022, 2021.
- 499 [56] Alireza Zaeemzadeh, Niccoló Bisagno, Zeno Sambugaro, Nicola Conci, Nazanin Rahnavard,  
500 and Mubarak Shah. Out-of-distribution detection using union of 1-dimensional subspaces. In  
501 *CVPR*, 2021.
- 502 [57] Joost Van Amersfoort, Lewis Smith, Yee Whye Teh, and Yarin Gal. Uncertainty estimation  
503 using a single deep deterministic neural network. In *ICML*, 2020.
- 504 [58] Sachin Vernekar, Ashish Gaurav, Vahdat Abdelzad, Taylor Denouden, Rick Salay, and Krzysztof  
505 Czarnecki. Out-of-distribution detection in classifiers via generation. In *NeurIPS Workshop*,  
506 2019.
- 507 [59] Alex Krizhevsky and Geoff Hinton. Convolutional deep belief networks on cifar-10. *Technical*  
508 *report, Citeseer*, 2009.
- 509 [60] Fisher Yu, Yinda Zhang, Shuran Song, Ari Seff, and Jianxiong Xiao. LSUN: construction of a  
510 large-scale image dataset using deep learning with humans in the loop. *CoRR*, abs/1506.03365,  
511 2015.
- 512 [61] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*,  
513 2015.

- 514 [62] Arthur Gretton, Karsten M. Borgwardt, Malte J. Rasch, Bernhard Schölkopf, and Alexander J.  
515 Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 2012.
- 516 [63] Itay Safran and Ohad Shamir. Depth-width tradeoffs in approximating natural functions with  
517 neural networks. In *ICML*, 2017.
- 518 [64] Allan Pinkus. Approximation theory of the mlp model in neural networks. *Acta numerica*,  
519 8:143–195, 1999.
- 520 [65] Peter L Bartlett and Wolfgang Maass. Vapnik-chervonenkis dimension of neural nets. *The  
521 handbook of brain theory and neural networks*, 2003.
- 522 [66] Ryuichi Kiryo, Gang Niu, Marthinus Christoffel du Plessis, and Masashi Sugiyama. Positive-  
523 unlabeled learning with non-negative risk estimator. In *NeurIPS*, 2017.
- 524 [67] Takashi Ishida, Gang Niu, and Masashi Sugiyama. Binary classification from positive-  
525 confidence data. In *NeurIPS*, 2018.
- 526 [68] Shuo Chen, Gang Niu, Chen Gong, Jun Li, Jian Yang, and Masashi Sugiyama. Large-margin  
527 contrastive learning with distance polarization regularizer. In *ICML*, 2021.

## 528 Checklist

- 529 1. For all authors...
- 530 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
531 contributions and scope? [Yes]
- 532 (b) Did you describe the limitations of your work? [Yes] See Appendix B
- 533 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See  
534 Appendix B
- 535 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
536 them? [Yes]
- 537 2. If you are including theoretical results...
- 538 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 539 (b) Did you include complete proofs of all theoretical results? [Yes]
- 540 3. If you ran experiments...
- 541 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
542 mental results (either in the supplemental material or as a URL)? [N/A]
- 543 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
544 were chosen)? [N/A]
- 545 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
546 ments multiple times)? [N/A]
- 547 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
548 of GPUs, internal cluster, or cloud provider)? [N/A]
- 549 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 550 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 551 (b) Did you mention the license of the assets? [N/A]
- 552 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 553
- 554 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
555 using/curating? [N/A]
- 556 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
557 information or offensive content? [N/A]
- 558 5. If you used crowdsourcing or conducted research with human subjects...
- 559 (a) Did you include the full text of instructions given to participants and screenshots, if  
560 applicable? [N/A]

561  
562  
563  
564

- (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]